

Appendix for Disentangling Disentanglement in Variational Autoencoders

Emile Mathieu*, Tom Rainforth*, N. Siddharth*, Yee Whye Teh

A. Proofs for Disentangling the β -VAE

Theorem 1. *The β -VAE target $\mathcal{L}_\beta(\mathbf{x})$ can be interpreted in terms of the standard ELBO, $\mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi)$, for an adjusted target $\pi_{\theta, \beta}(\mathbf{x}, \mathbf{z}) \triangleq p_\theta(\mathbf{x} | \mathbf{z})f_\beta(\mathbf{z})$ with annealed prior $f_\beta(\mathbf{z}) \triangleq p(\mathbf{z})^\beta / F_\beta$ as*

$$\mathcal{L}_\beta(\mathbf{x}) = \mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi) + (\beta - 1)H_{q_\phi} + \log F_\beta \quad (3)$$

where $F_\beta \triangleq \int_{\mathbf{z}} p(\mathbf{z})^\beta d\mathbf{z}$ is constant given β , and H_{q_ϕ} is the entropy of $q_\phi(\mathbf{z} | \mathbf{x})$.

Proof. Starting with (2), we have

$$\begin{aligned} \mathcal{L}_\beta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})}[\log p_\theta(\mathbf{x} | \mathbf{z})] + \beta H_{q_\phi} \\ &\quad + \beta \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})}[\log p(\mathbf{z})] \\ &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})}[\log p_\theta(\mathbf{x} | \mathbf{z})] + (\beta - 1)H_{q_\phi} + H_{q_\phi} \\ &\quad + \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})}[\log p(\mathbf{z})^\beta - \log F_\beta] + \log F_\beta \\ &= \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})}[\log p_\theta(\mathbf{x} | \mathbf{z})] + (\beta - 1)H_{q_\phi} \\ &\quad - \text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) \| f_\beta(\mathbf{z})) + \log F_\beta \\ &= \mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi) + (\beta - 1)H_{q_\phi} + \log F_\beta \end{aligned}$$

as required. \square

Corollary 1. *If $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \Sigma)$ and $q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_\phi(\mathbf{x}), S_\phi(\mathbf{x}))$, then,*

$$\mathcal{L}_\beta(\mathbf{x}; \theta, \phi) = \mathcal{L}(\mathbf{x}; \theta', \phi') + \frac{(\beta - 1)}{2} \log |S_{\phi'}(\mathbf{x})| + c \quad (4)$$

where θ' and ϕ' represent rescaled networks such that

$$\begin{aligned} p_{\theta'}(\mathbf{x} | \mathbf{z}) &= p_\theta(\mathbf{x} | \mathbf{z} / \sqrt{\beta}), \\ q_{\phi'}(\mathbf{z} | \mathbf{x}) &= \mathcal{N}(\mathbf{z}; \mu_{\phi'}(\mathbf{x}), S_{\phi'}(\mathbf{x})), \\ \mu_{\phi'}(\mathbf{x}) &= \sqrt{\beta} \mu_\phi(\mathbf{x}), \quad S_{\phi'}(\mathbf{x}) = \beta S_\phi(\mathbf{x}), \end{aligned}$$

and $c \triangleq \frac{D(\beta-1)}{2} \left(1 + \log \frac{2\pi}{\beta}\right) + \log F_\beta$ is a constant, with D denoting the dimensionality of \mathbf{z} .

Proof. We start by noting that

$$\begin{aligned} \pi_{\theta, \beta}(\mathbf{x}) &= \mathbb{E}_{f_\beta(\mathbf{z})}[p_\theta(\mathbf{x} | \mathbf{z})] = \mathbb{E}_{p(\mathbf{z})}[p_\theta(\mathbf{x} | \mathbf{z} / \sqrt{\beta})] \\ &= \mathbb{E}_{p(\mathbf{z})}[p_{\theta'}(\mathbf{x} | \mathbf{z})] = p_{\theta'}(\mathbf{x}) \end{aligned}$$

Now considering an alternate form of $\mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi)$ in (3),

$$\begin{aligned} \mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi) &= \log \pi_{\theta, \beta}(\mathbf{x}) - \text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) \| \pi_{\theta, \beta}(\mathbf{z} | \mathbf{x})) \\ &= \log p_{\theta'}(\mathbf{x}) - \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{q_\phi(\mathbf{z} | \mathbf{x}) p_{\theta'}(\mathbf{x})}{p_\theta(\mathbf{x} | \mathbf{z}) f_\beta(\mathbf{z})} \right) \right] \end{aligned}$$

$$\begin{aligned} &= \log p_{\theta'}(\mathbf{x}) \\ &\quad - \mathbb{E}_{q_{\phi'}(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{q_\phi(\mathbf{z} / \sqrt{\beta} | \mathbf{x}) p_{\theta'}(\mathbf{x})}{p_\theta(\mathbf{x} | \mathbf{z} / \sqrt{\beta}) f_\beta(\mathbf{z} / \sqrt{\beta})} \right) \right]. \quad (8) \end{aligned}$$

We first simplify $f_\beta(\mathbf{z} / \sqrt{\beta})$ as

$$\begin{aligned} f_\beta(\mathbf{z} / \sqrt{\beta}) &= \frac{1}{\sqrt{2\pi}^{|\Sigma|} |\Sigma|^{|\beta|}} \exp\left(-\frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z}\right) \\ &= p(\mathbf{z}) \beta^{(D/2)}. \end{aligned}$$

Further, denoting $\mathbf{z}_\ddagger = \mathbf{z} - \sqrt{\beta} \mu_{\phi'}(\mathbf{x})$, and $\mathbf{z}_\ddagger = \mathbf{z}_\ddagger / \sqrt{\beta} = \mathbf{z} / \sqrt{\beta} - \mu_{\phi'}(\mathbf{x})$, we have

$$\begin{aligned} q_{\phi'}(\mathbf{z} | \mathbf{x}) &= \frac{1}{\sqrt{2\pi}^{|\Sigma_\phi(\mathbf{x})|} |\Sigma_\phi(\mathbf{x})|^{|\beta|}} \exp\left(-\frac{1}{2\beta} \mathbf{z}_\ddagger^T S_\phi(\mathbf{x})^{-1} \mathbf{z}_\ddagger\right), \\ q_\phi\left(\frac{\mathbf{z}}{\sqrt{\beta}} | \mathbf{x}\right) &= \frac{1}{\sqrt{2\pi}^{|\Sigma_\phi(\mathbf{x})|} |\Sigma_\phi(\mathbf{x})|^{|\beta|}} \exp\left(-\frac{1}{2} \mathbf{z}_\ddagger^T S_\phi(\mathbf{x})^{-1} \mathbf{z}_\ddagger\right) \\ &\quad \text{giving } q_\phi\left(\mathbf{z} / \sqrt{\beta} | \mathbf{x}\right) = q_{\phi'}(\mathbf{z} | \mathbf{x}) \beta^{(D/2)}. \end{aligned}$$

Plugging these back in to (8) while remembering $p_\theta(\mathbf{x} | \mathbf{z} / \sqrt{\beta}) = p_{\theta'}(\mathbf{x} | \mathbf{z})$, we have

$$\begin{aligned} \mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi) &= \log p_{\theta'}(\mathbf{x}) - \mathbb{E}_{q_{\phi'}(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{q_{\phi'}(\mathbf{z} | \mathbf{x}) p_{\theta'}(\mathbf{x})}{p_{\theta'}(\mathbf{x} | \mathbf{z}) p(\mathbf{z})} \right) \right] \\ &= \mathcal{L}(\mathbf{x}; \theta, \phi), \end{aligned}$$

showing that the ELBOs for the two setups are the same. For the entropy term, we note that

$$\begin{aligned} H_{q_\phi} &= \frac{D}{2} (1 + \log 2\pi) + \frac{1}{2} \log |S_\phi(\mathbf{x})| \\ &= \frac{D}{2} \left(1 + \log \frac{2\pi}{\beta}\right) + \frac{1}{2} \log |S_{\phi'}(\mathbf{x})|. \end{aligned}$$

Finally substituting for H_{q_ϕ} and $\mathcal{L}(\mathbf{x}; \pi_{\theta, \beta}, q_\phi)$ in (3) gives the desired result. \square

Corollary 2. *Let $[\theta', \phi'] = g_\beta([\theta, \phi])$ represent the transformation required to produce the rescaled networks in Corollary 1. If $0 < |\det \nabla_{\theta, \phi} g([\theta, \phi])| < \infty \forall [\theta, \phi]$, then*

$$\nabla_{\theta, \phi} \mathcal{L}_\beta(\mathbf{x}; \theta, \phi) = \mathbf{0} \Leftrightarrow \nabla_{\theta', \phi'} \mathcal{L}_{H, \beta}(\mathbf{x}; \theta', \phi') = \mathbf{0}.$$

Thus $[\theta^*, \phi^*]$ being a stationary point of $\frac{1}{n} \sum_{i=1}^n \mathcal{L}_\beta(\mathbf{x}_i; \theta, \phi)$ indicates that $g_\beta([\theta^*, \phi^*])$ is a stationary point of $\frac{1}{n} \sum_{i=1}^n \mathcal{L}_{H, \beta}(\mathbf{x}_i; \theta', \phi')$ and vice-versa.

Proof. Starting from Corollary 1 we have that

$$\begin{aligned} \nabla_{\theta, \phi} \mathcal{L}_\beta(\mathbf{x}; \theta, \phi) &= \nabla_{\theta, \phi} \mathcal{L}_{H, \beta}(\mathbf{x}; \theta', \phi') \\ &= (\nabla_{\theta, \phi} g_\beta([\theta, \phi])) \nabla_{\theta', \phi'} \mathcal{L}_{H, \beta}(\mathbf{x}; \theta', \phi'), \end{aligned}$$

so $\nabla_{\theta', \phi'} \mathcal{L}_{H, \beta}(\mathbf{x}; \theta', \phi') = \mathbf{0} \implies \nabla_{\theta, \phi} \mathcal{L}_{\beta}(\mathbf{x}; \theta, \phi) = \mathbf{0}$ given our assumption that $|\det \nabla_{\theta, \phi} g([\theta, \phi])| < \infty \forall [\theta, \phi]$. Further, as $0 < |\det \nabla_{\theta, \phi} g([\theta, \phi])| \forall [\theta, \phi]$, $(\nabla_{\theta, \phi} g([\theta, \phi]))^{-1}$ exists and has a finite determinant, so $\nabla_{\theta, \phi} \mathcal{L}_{\beta}(\mathbf{x}; \theta, \phi) = \mathbf{0}$ also implies $\nabla_{\theta', \phi'} \mathcal{L}_{H, \beta}(\mathbf{x}; \theta', \phi') = \mathbf{0}$. \square

Theorem 2. If $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \sigma I)$ and $q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\phi}(\mathbf{x}), S_{\phi}(\mathbf{x}))$, then for all rotation matrices R ,

$$\mathcal{L}_{\beta}(\mathbf{x}; \theta, \phi) = \mathcal{L}_{\beta}(\mathbf{x}; \theta^{\dagger}(R), \phi^{\dagger}(R)) \quad (6)$$

where $\theta^{\dagger}(R)$ and $\phi^{\dagger}(R)$ are transformed networks such that

$$\begin{aligned} p_{\theta^{\dagger}}(\mathbf{x} | \mathbf{z}) &= p_{\theta}(\mathbf{x} | R^T \mathbf{z}), \\ q_{\phi^{\dagger}}(\mathbf{z} | \mathbf{x}) &= \mathcal{N}(\mathbf{z}; R \mu_{\phi}(\mathbf{x}), R S_{\phi}(\mathbf{x}) R^T). \end{aligned}$$

Proof. If $\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})$ and $\mathbf{y} = R\mathbf{z}$ then, by Petersen et al. (§8.1.4 2008)), we have

$$\mathbf{y} \sim \mathcal{N}(\mathbf{y}; R \mu_{\phi}(\mathbf{x}), R S_{\phi}(\mathbf{x}) R^T).$$

Consequently, the changes made by the transformed networks cancel to give the same reconstruction error as

$$\begin{aligned} \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] &= \mathbb{E}_{q_{\phi^{\dagger}}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | R^T \mathbf{z})] \\ &= \mathbb{E}_{q_{\phi^{\dagger}}(\mathbf{z} | \mathbf{x})} [\log p_{\theta^{\dagger}}(\mathbf{x} | \mathbf{z})]. \end{aligned}$$

Furthermore, the KL divergence between $q_{\phi}(\mathbf{z} | \mathbf{x})$ and $p_{\theta}(\mathbf{z})$ is invariant to rotation, because of the rotational symmetry of the latter, such that $\text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})) = \text{KL}(q_{\phi^{\dagger}}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}))$. The result now follows from noting that the two terms of the β -VAE are equal under rotation. \square

B. Experimental Details

Disentanglement - 2d-shapes: The experiments from Section 6 on the impact of the prior in terms of disentanglement are conducted on the **2D Shapes** (Matthey et al., 2017) dataset, comprising of 737,280 binary 64 x 64 images of 2D shapes with ground truth factors [number of values]: shape [3], scale [6], orientation [40], x-position [32], y-position [32]. We use a convolutional neural network for the encoder and a deconvolutional neural network for the decoder, whose architectures are described in Table 1a. We use $[0, 1]$ normalised data as targets for the mean of a Bernoulli distribution and negative cross-entropy for $\log p(\mathbf{x} | \mathbf{z})$. We rely on the Adam optimiser (Kingma and Ba, 2015; Reddi et al., 2018) with learning rate $1e^{-4}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$, to optimise the β -VAE objective from (3).

For $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \text{diag}(\sigma))$, experiments were run with a batch size of 64 and for 20 epochs. For $p(\mathbf{z}) = \prod_d \text{STUDENT-T}(z_d; \nu)$, experiments were run with a batch size of 256 and for 40 epochs. In Figure 2, the *PCA initialised anisotropic* prior is initialised so that its standard deviations are set to be the first D singular values of the data. These are then mapped through a softmax function to ensure that the β regularisation coefficient is not imple-

Encoder	Decoder
Input 64 x 64 binary image	Input $\in \mathbb{R}^{10}$
4x4 conv. 32 stride 2 & ReLU	FC. 128 ReLU
4x4 conv. 32 stride 2 & ReLU	FC. 4x4 x 64 ReLU
4x4 conv. 64 stride 2 & ReLU	4x4 upconv. 64 stride 2 & ReLU
4x4 conv. 64 stride 2 & ReLU	4x4 upconv. 64 stride 2 & ReLU
FC. 128	4x4 upconv. 32 stride 2 & ReLU
FC. 2x10	4x4 upconv. 1. stride 2

(a) 2D-shapes dataset.

Encoder	Decoder
Input $\in \mathbb{R}^2$	Input $\in \mathbb{R}^2$
FC. 100. & ReLU	FC. 100 & ReLU
FC. 2x2	FC. 2x2

(b) Pinwheel dataset.

Encoder
Input 32 x 32 x 1 channel image
4x4 conv. 32 stride 2 & BatchNorm2d & LeakyReLU(.2)
4x4 conv. 64 stride 2 & BatchNorm2d & LeakyReLU(.2)
4x4 conv. 128 stride 2 & BatchNorm2d & LeakyReLU(.2)
4x4 conv. 50, 4x4 conv. 50
Decoder
Input $\in \mathbb{R}^{50}$
4x4 upconv. 128 stride 1 pad 0 & BatchNorm2d & ReLU
4x4 upconv. 64 stride 2 pad 1 & BatchNorm2d & ReLU
4x4 upconv. 32 stride 2 pad 1 & BatchNorm2d & ReLU
4x4 upconv. 1 stride 2 pad 1

(c) Fashion-MNIST dataset.

Table 1. Encoder and decoder architectures.

itly scaled compared to the isotropic case. For the *learned anisotropic* priors, standard deviations are first initialised as just described, and then learned along with the model through a log-variance parametrisation.

We rely on the metric presented in §4 and Appendix B of Kim and Mnih (2018) as a measure of axis-alignment of the latent encodings with respect to the true (known) generative factors. Confidence intervals in Figure 2 were computed via the assumption of normally distributed samples with unknown mean and variance, with 100 runs of each model.

Clustering - Pinwheel We generated spiral cluster data¹, with $n = 400$ observations, clustered in 4 spirals, with radial and tangential standard deviations respectively of 0.1 and 0.30, and a rate of 0.25. We use fully-connected neural networks for both the encoder and decoder, whose architectures are described in Table 1b. We minimise the objective from (7), with \mathbb{D} chosen to be the inclusive KL and $q_{\phi}(\mathbf{z})$ approximated by the aggregate encoding of the full dataset:

$$\begin{aligned} \mathbb{D}(q_{\phi}(\mathbf{z}), p(\mathbf{z})) &= \text{KL}(p(\mathbf{z}) \| q_{\phi}(\mathbf{z})) \\ &= \mathbb{E}_{p(\mathbf{z})} [\log(p(\mathbf{z})) - \log(\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})} [q_{\phi}(\mathbf{z} | \mathbf{x})])] \end{aligned}$$

¹<http://hips.seas.harvard.edu/content/synthetic-pinwheel-data-matlab>.

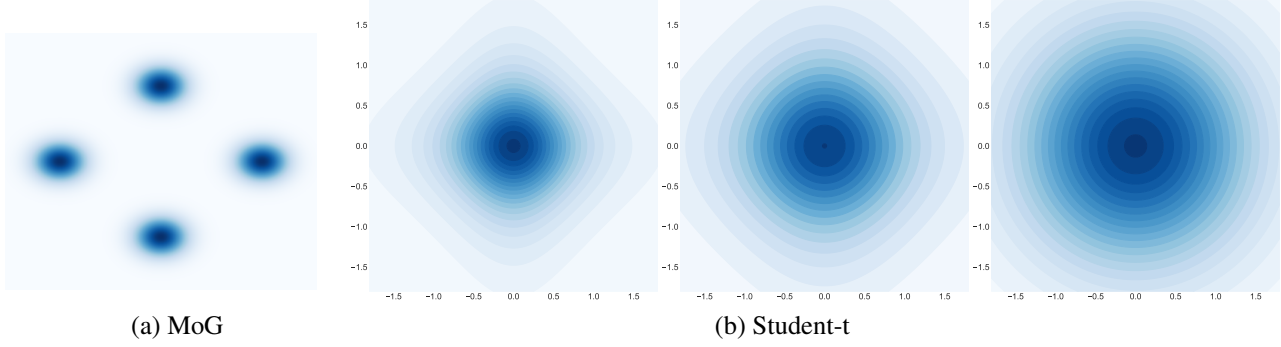


Figure 6. (a) PDF of Gaussian mixture model prior $p(\mathbf{z})$, as per (9). (b) PDF for a 2-dimensional factored Student-t distributions p_ν with degree of freedom $\nu = \{3, 5, 100\}$ (left to right). Note that $p_\nu(\mathbf{z}) \rightarrow \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ as $\nu \rightarrow \infty$.

$$\approx \sum_{j=1}^B \left(\log p(\mathbf{z}_j) - \log \left(\sum_{i=1}^n q_\phi(\mathbf{z}_j | \mathbf{x}_i) \right) \right)$$

with $\mathbf{z}_j \sim p(\mathbf{z})$. A Gaussian likelihood is used for the encoder. We trained the model for 500 epochs using the Adam optimiser (Kingma and Ba, 2015; Reddi et al., 2018), with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ and a learning rate of $1e^{-3}$. The batch size is set to $B = n$.

The Gaussian mixture prior (c.f. Figure 6(a)) is defined as

$$p(\mathbf{z}) = \sum_{c=1}^C \pi^c \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}^c, \boldsymbol{\Sigma}^c) \\ = \sum_{c=1}^C \pi^c \prod_{d=1}^D \mathcal{N}(z_d | \mu_d^c, \sigma_d^c) \quad (9)$$

with $D = 2$, $C = 4$, $\boldsymbol{\Sigma}^c = 0.03I_D$, $\pi^c = 1/C$, and $\mu_d^c \in \{0, 1\}$.

Sparsity - Fashion-MNIST The experiments from Section 6 on the latent representation’s sparsity are conducted on the **Fashion-MNIST** (Xiao et al., 2017) dataset, comprising of 70,000 greyscale images resized to 32×32 .

To enforce sparsity, we relied on a prior defined as a factored univariate mixture of a standard and low variance normal distributions:

$$p(z_d) = \prod_d (1 - \gamma) \mathcal{N}(z_d; 0, 1) + \gamma \mathcal{N}(z_d; 0, \sigma_0^2)$$

with $\sigma_0^2 = 0.05$. The weight, γ , of the low-variance component indicates how likely samples are to come from that component, hence to be *off*.

We minimised the objective from (7), with $\mathbb{D}(q_\phi(\mathbf{z}), p(\mathbf{z}))$ taken to be a dimension-wise MMD with a sum of *Cauchy* kernels on each dimension. Equivalently, we can think of this as calculating a single MMD using the single kernel

$$k(\mathbf{x}, \mathbf{y}) = \sum_{d=1}^D \sum_{\ell} \frac{\sigma_\ell}{\sigma_{\ell=1} + (x_d - y_d)^2}. \quad (10)$$

where $\sigma_\ell \in \{0.2, 0.4, 1, 2, 4, 10\}$ are a set of length scales.

This dimension-wise kernel only enforces a congruence between the marginal distributions of \mathbf{x} and \mathbf{y} and so, strictly speaking, its MMD does not constitute a valid divergence metric in the sense that we can have $\mathbb{D}(q_\phi(\mathbf{z}), p(\mathbf{z})) = 0$ when $q_\phi(\mathbf{z})$ and $p(\mathbf{z})$ are not identical distributions: it only requires their marginals to match to get zero divergence.

The reasons we chose this approach are twofold. Firstly, we found that conventional kernels based on the Euclidean distance between encodings produced gradients with insurmountably high variances, meaning that effectively minimizing the divergence to get $q_\phi(\mathbf{z})$ and $p(\mathbf{z})$ to match was not possible, even for very large batch sizes and $\alpha \rightarrow \infty$.

Secondly, though just matching the marginal distributions is not sufficient to ensure sparsity—as one could have some points with all dimensions close to the origin and some with all dimensions far away—a combination of the need to achieve good reconstructions and noise in the encoder process should prevent this from occurring. In short, provided the noise from the encoder is properly regulated, there is little information that can be stored in latent dimensions near the origin because of the high level of overlap forced in this region. Therefore, for a datapoint to be effectively encoded, it must have at least some of its latent dimensions outside of this region. Coupled with the need for most of the latent values to be near the origin to match the marginal distributions, this, in turn, enforces a sparse representation. Consequently, the loss in sparsity performance relative to using a hypothetical kernel that is both universal and has stable gradient estimates should only be relatively small, as is borne out in our empirical results. This may, however, be why we see a slight drop in sparsity performance for very large values of α .

We use a convolutional neural network for the encoder and a deconvolutional neural network for the decoder, whose architectures come from the DCGAN model (Radford et al., 2016) and are described in Table 1c. We use $[0, 1]$ normalised data as targets for the mean of a Laplace distribution with fixed scaling of 0.1. We rely on the Adam optimiser

with learning rate $5e^{-4}$, $\beta_1 = 0.5$, and $\beta_2 = 0.999$. The model is then trained (on the training set) for 80 epochs with a batch-size of 500.

As an extrinsic measure of *sparsity*, we use the *Hoyer* metric (Hurley and Rickard, 2008), defined for $\mathbf{y} \in \mathbb{R}^d$ by

$$\text{Hoyer}(\mathbf{y}) = \frac{\sqrt{d} - \|\mathbf{y}\|_1 / \|\mathbf{y}\|_2}{\sqrt{d} - 1} \in [0, 1],$$

yielding 0 for a fully dense vector and 1 for a fully sparse vector. We additionally normalise each dimension to have a standard deviation of 1 under its aggregate distribution, i.e. we use $\bar{z}_d = z_d / \sigma(z_d)$ where $\sigma(z_d)$ is the standard deviation of dimension d of the latent encoding taken over the dataset. Overall sparsity is computed by averaging over the dataset as $\text{Sparsity} = 1/n \sum_i \text{Hoyer}(\bar{z}_i)$.

As discussed in the main text, we use a trained model with $\alpha = 1000$, $\beta = 1$, and $\gamma = 0.8$ to perform a qualitative analysis of sparsity using the *Fashion-MNIST* dataset. Figure 7 shows the per-class average embedding magnitude for this model, a subset of which was shown in the main text. As can be seen clearly, the different classes utilise predominantly different subsets of dimensions to encode the image data, as one might expect for sparse representations.

C. Posterior regularisation

The aggregate posterior regulariser $\mathbb{D}(q(z), p(z))$ is a little more subtle to analyse than the entropy regulariser as it involves both the choice of divergence and potential difficulties in estimating that divergence. One possible choice is the exclusive Kullback-Leibler divergence $\text{KL}(q(z) \parallel p(z))$, as previously used (without additional entropy regularisation) by (Dilokthanakul et al., 2019; Esmaeili et al., 2019), but also implicitly by (Chen et al., 2018), through the use of a total correlation (TC) term. We now highlight a shortfall with this choice of divergence due to difficulties in its empirical estimation.

In short, the approaches used to estimate the $\text{H}[q(z)]$ (noting that $\text{KL}(q(z) \parallel p(z)) = -\text{H}[q(z)] - \mathbb{E}_{q(z)}[\log p(z)]$, where the latter term can be estimated reliably by a simple Monte Carlo estimate) can exhibit very large biases unless very large batch sizes are used, resulting in quite different effects from what was intended. In fact, our results suggest they will exhibit behaviour similar to the β -VAE if the batch size is too small. These biases arise from the effects of nesting estimators (Rainforth et al., 2018a), where the variance in the nested (inner) estimator for $q(z)$ induces a bias in the overall estimator. Specifically, for any random variable \hat{Z} ,

$$\mathbb{E}[\log(\hat{Z})] = \log(\mathbb{E}[\hat{Z}]) - \frac{\text{Var}[\hat{Z}]}{2Z^2} + O(\varepsilon)$$

where $O(\varepsilon)$ represents higher-order moments that get dominated asymptotically if \hat{Z} is a Monte-Carlo estimator (see

Proposition 1c in Maddison et al. (2017), Theorem 1 in Rainforth et al. (2018b), or Theorem 3 in Domke and Sheldon (2018)). In this setting, $\hat{Z} = \hat{q}(z)$ is the estimate used for $q(z)$. We thus see that if the variance of $\hat{q}(z)$ is large, this will induce a significant bias in our KL estimator.

To make things precise, we consider the estimator used for $\text{H}[q(z)]$ in Chen et al. (2018); Dilokthanakul et al. (2019); Esmaeili et al. (2019)

$$\text{H}[q(z)] \approx \hat{\text{H}} \triangleq -\frac{1}{B} \sum_{b=1}^B \log \hat{q}(z_b), \text{ where} \quad (11a)$$

$$\hat{q}(z_b) = \frac{q_\phi(z_b | \mathbf{x}_b)}{n} + \frac{n-1}{n(B-1)} \sum_{b' \neq b} q_\phi(z_b | \mathbf{x}'_{b'}), \quad (11b)$$

$z_b \sim q_\phi(z | \mathbf{x}_b)$, and $\{\mathbf{x}_1, \dots, \mathbf{x}_B\}$ is the mini-batch of data used for the current iteration for dataset size n . Esmaeili et al. (2019) correctly show that $\mathbb{E}[\hat{q}(z_b)] = \tilde{q}(z_b)$, with the first term of (11b) comprising an exact term in $\tilde{q}(z_b)$ and the second term of (11b) being an unbiased Monte-Carlo estimate for the remaining terms in $\tilde{q}(z_b)$.

To examine the practical behaviour of this estimator when $B \ll n$, we first note that the second term of (11b) is, in practice, usually very small and dominated by the first term. This is borne out empirically in our own experiments, and also noted in Kim and Mnih (2018). To see why this is the case, consider that given encodings of two independent data points, it is highly unlikely that the two encoding distributions will have any notable overlap (e.g. for a Gaussian encoder, the means will most likely be very many standard deviations apart), presuming a sensible latent space is being learned. Consequently, even though this second term is unbiased and may have an expectation comparable or even larger than the first, it is heavily skewed—it is usually negligible, but occasionally large in the rare instances where there is substantial overlap between encodings.

Let the second term of (11b) be T_2 and the event that this is significant be E_S , such that $\mathbb{E}[T_2 | \neg E_S] \approx 0$. As explained above, $\mathbb{P}(E_S) \ll 1$ typically. We now have

$$\begin{aligned} \mathbb{E}[\hat{\text{H}}] &= \mathbb{P}(E_S) \mathbb{E}[\hat{\text{H}} | E_S] + (1 - \mathbb{P}(E_S)) \mathbb{E}[\hat{\text{H}} | \neg E_S] \\ &= \mathbb{P}(E_S) \mathbb{E}[\hat{\text{H}} | E_S] + (1 - \mathbb{P}(E_S)) \\ &\quad \cdot (\log n - \frac{1}{B} \sum_{b=1}^B \mathbb{E}[\log q_\phi(z_b | \mathbf{x}_b) | \neg E_S] - \mathbb{E}[T_2 | \neg E_S]) \\ &= \mathbb{P}(E_S) \mathbb{E}[\hat{\text{H}} | E_S] + (1 - \mathbb{P}(E_S)) \\ &\quad \cdot (\log n - \mathbb{E}[\log q_\phi(z_1 | \mathbf{x}_1) | \neg E_S] - \mathbb{E}[T_2 | \neg E_S]) \\ &\approx \mathbb{P}(E_S) \mathbb{E}[\hat{\text{H}} | E_S] \\ &\quad + (1 - \mathbb{P}(E_S)) (\log n - \mathbb{E}[\log q_\phi(z_1 | \mathbf{x}_1)]) \end{aligned}$$

where the approximation relies firstly on our previous assumption that $\mathbb{E}[T_2 | \neg E_S] \approx 0$ and also that $\mathbb{E}[\log q_\phi(z_1 | \mathbf{x}_1) | \neg E_S] \approx \mathbb{E}[\log q_\phi(z_1 | \mathbf{x}_1)]$. This second

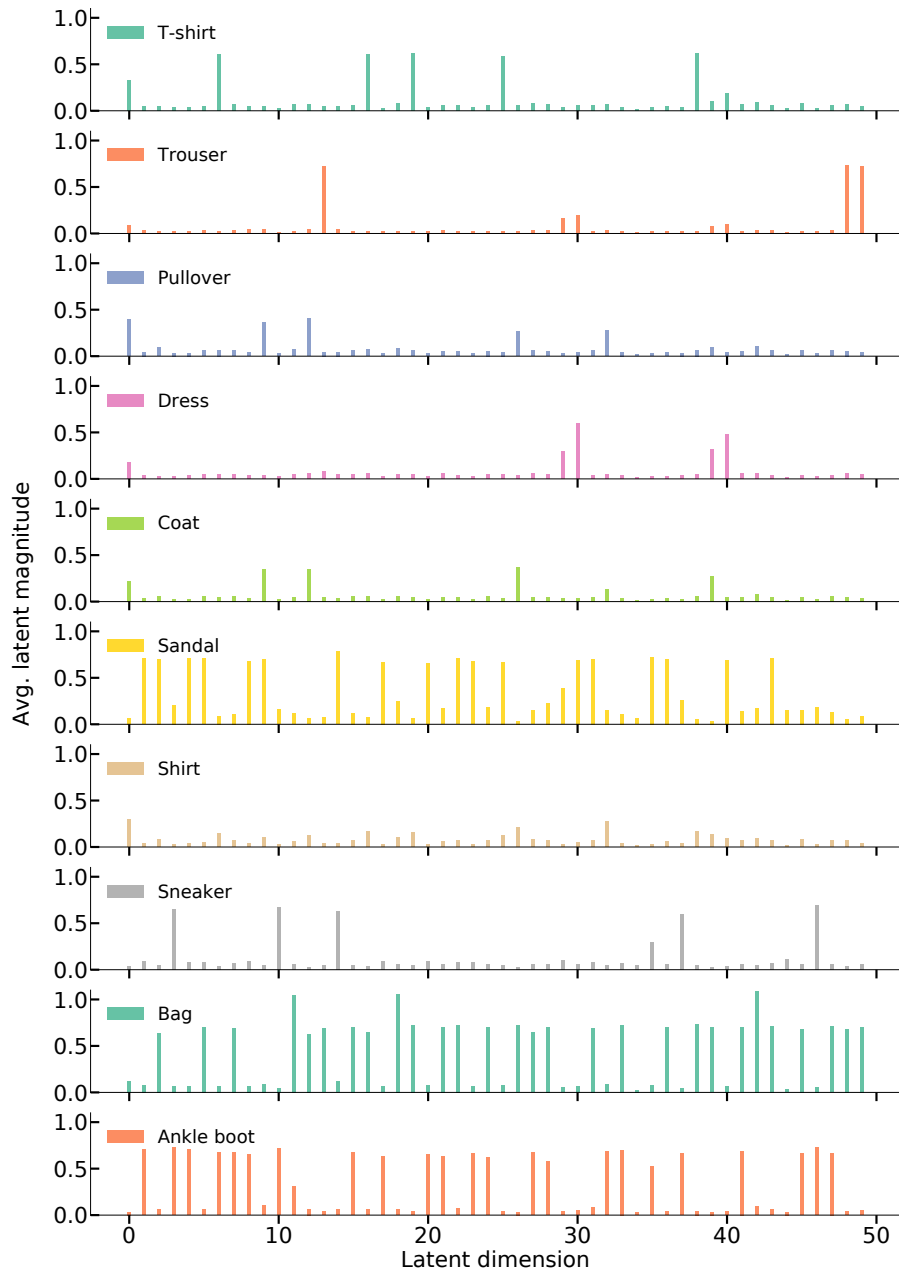


Figure 7. Average encoding magnitude over data for each classes in *Fashion-MNIST*.

assumption will also generally hold in practice, firstly because the occurrence of E_S is dominated by whether two similar datapoints are drawn (rather than by the value of \mathbf{x}_1) and secondly because $\mathbb{P}(E_S) \ll 1$ implies that

$$\begin{aligned} & \mathbb{E}[\log q_\phi(\mathbf{z}_1 | \mathbf{x}_1)] \\ &= (1 - \mathbb{P}(E_S)) \mathbb{E}[\log q_\phi(\mathbf{z}_1 | \mathbf{x}_1) | \neg E_S] \\ & \quad + \mathbb{P}(E_S) \mathbb{E}[\log q_\phi(\mathbf{z}_1 | \mathbf{x}_1) | E_S] \\ & \approx \mathbb{E}[\log q_\phi(\mathbf{z}_1 | \mathbf{x}_1) | \neg E_S]. \end{aligned}$$

Characterising $\mathbb{E}[\hat{\mathbf{H}} | E_S]$ precisely is a little more challenging, but it can safely be assumed to be smaller than $\mathbb{E}[\log q_\phi(\mathbf{z}_1 | \mathbf{x}_1)]$, which is approximately what would result from all the \mathbf{x}'_b being the same as \mathbf{x}_b . We thus see that even when the event E_S does occur, the resulting estimates will still, at most, be on a comparable scale to when it does not. Consequently, whenever E_S is rare, the $(1 - \mathbb{P}(E_S)) \mathbb{E}[\hat{\mathbf{H}} | \neg E_S]$ term will dominate and we thus have

$$\begin{aligned} \mathbb{E}[\hat{\mathbf{H}}] & \approx \log n - \mathbb{E}[\log q_\phi(\mathbf{z}_1 | \mathbf{x}_1)] \\ & = \log n + \mathbb{E}_{p(\mathbf{x})}[\mathbf{H}[q_\phi(\mathbf{z} | \mathbf{x})]]. \end{aligned}$$

We now see that the estimator mimics the β -VAE regularisation up to a constant factor $\log n$, as adding the $\mathbb{E}_{q(\mathbf{z})}[\log p(\mathbf{z})]$ back in gives

$$\begin{aligned} & - \mathbb{E}[\hat{\mathbf{H}}] - \mathbb{E}_{q(\mathbf{z})}[\log p(\mathbf{z})] \\ & \approx \mathbb{E}_{p(\mathbf{x})}[\text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}))] - \log n. \end{aligned}$$

We should thus expect to empirically see training with this estimator as a regulariser to behave similarly to the β -VAE with the same regularisation term whenever $B \ll n$. Note that the $\log n$ constant factor will not impact the gradients, but does mean that it is possible, even likely, that negative estimates for $\hat{\mathbf{K}}\mathbf{L}$ will be generated, even though we know the true value is positive.

Overcoming the problem can, at least to a certain degree, be overcome by using very large batch sizes B , at an inevitable computational and memory cost. However, the problem is potentially exacerbated in higher dimensional latent spaces and larger datasets, for which one would typically expect the typical overlap of datapoints to decrease.

C.1. Other Divergences

As discussed in the main paper, $\text{KL}(q(\mathbf{z}) \| p(\mathbf{z}))$ is far from the only aggregate posterior regulariser one might use. Though we do not analyse them formally, we expect many alternative divergence-estimator pairs to suffer from similar issues. For example, using Monte Carlo estimators with the inclusive Kullback-Leibler divergence $\text{KL}(p(\mathbf{z}) \| q(\mathbf{z}))$ or the sliced Wasserstein distance (Kolouri et al., 2019) both result in nested expectations analogously to $\text{KL}(q(\mathbf{z}) \| p(\mathbf{z}))$, and are therefore likely to similarly induce substantial bias without using large batch sizes.

Interestingly, however, MMD and generative adversarial network (GAN) regularisers of the form discussed in (Tolstikhin et al., 2018) do not result in nested expectations and therefore are necessarily not prone to the same issues: they produce unbiased estimates of their respective objectives. Though we experienced practical issues in successfully implementing both of these—we found the signal-to-noise-ratio of the MMD gradient estimates to be very low, particularly in high dimensions, while we experienced training instabilities for the GAN regulariser—their apparent theoretical advantages may indicate that they are preferable approaches, particularly if these issues can be alleviated. The GAN-based approach to estimating the total correlation introduced by Kim and Mnih (2018) similarly allows a nested expectation to be avoided, at the cost of converting a conventional optimization into a minimax problem.

Given the failings of the available existing approaches, we believe that further investigation into divergence-estimator pairs for $\mathbb{D}(q(\mathbf{z}), p(\mathbf{z}))$ in VAEs is an important topic for future work that extends well beyond the context of this paper, or even the general aim of achieving decomposition. In particular, the need for congruence between the posterior (encoder), likelihood (decoder), and marginal likelihood (data distribution) for a generative model, means that ensuring $q(\mathbf{z})$ is close to $p(\mathbf{z})$ is a generally important endeavour for training VAEs. For example, mismatch between $q(\mathbf{z})$ and $p(\mathbf{z})$ will cause samples drawn from the learned generative model to mismatch the true data-generating distribution, regardless of the fidelity of our encoder and decoder.

D. Characterising Overlap

Reiterating the argument from the main text, although the mutual information $I(\mathbf{x}; \mathbf{z})$ between data and latents provides a perfectly serviceable characterisation of overlap in a number of cases, the two are not universally equivalent and we argue that it is overlap which is important in achieving useful representations. In particular, if the form of the encoding distribution is not fixed—as when employing normalising flows, for example— $I(\mathbf{x}; \mathbf{z})$ does not necessarily characterise overlap well.

Consider, for example, an encoding distribution that is a mixture between the prior and a uniform distribution on a tiny ϵ -ball around the mean encoding $\mu_\phi(\mathbf{x})$, i.e. $q_\phi(\mathbf{z} | \mathbf{x}) = \lambda \cdot \text{Uniform}(\|\mu_\phi(\mathbf{x}) - \mathbf{z}\|_2 < \epsilon) + (1 - \lambda) \cdot p(\mathbf{z})$. If the encoder and decoder are sufficiently flexible to learn arbitrary representations, one now could arrive at *any* value for mutual information simply by an appropriate choice of λ . However, enforcing structuring of the latent space will be effectively impossible due to the lack of any pressure (other than a potentially small amount from internal regularization in the encoder network itself) for similar encodings to correspond to similar datapoints; the overlap between any two

encodings is the same unless they are within ϵ of each other.

While this example is a bit contrived, it highlights a key feature of overlap that $I(\mathbf{x}; \mathbf{z})$ fails to capture: $I(\mathbf{x}; \mathbf{z})$ does not distinguish between large overlap with a small number of other datapoints and small overlap with a large number of other datapoints. This distinction is important because we are particularly interested in *how many* other datapoints one datapoint’s encoding overlaps with when imposing structure—the example setup fails because each datapoint has the same level of overlap with all the other datapoints.

Another feature that $I(\mathbf{x}; \mathbf{z})$ can fail to account for is a notion of *locality* in the latent space. Imagine a scenario where the encoding distributions are extremely multimodal with similar sized modes spread throughout the latent space, such as $q(\mathbf{z}|\mathbf{x}) = \sum_{i=1}^{1000} \mathcal{N}(\mathbf{z}; \mu_{\phi}(\mathbf{x}) + m_i, \sigma I)$ for some constant scalar σ , and vectors m_i . Again we can achieve almost any value for $I(\mathbf{x}; \mathbf{z})$ by adjusting σ , but it is difficult to impose meaningful structure regardless as each datapoint can be encoded to many different regions of the latent space.

References

- Ricky T. Q. Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, 2018.
- Nat Dilokthanakul, Nick Pawlowski, and Murray Shanahan. Explicit information placement on latent variables using auxiliary generative modelling task, 2019. URL <https://openreview.net/forum?id=H11-SjA5t7>.
- Justin Domke and Daniel Sheldon. Importance weighting and variational inference. In *Advances in Neural Information Processing Systems*, pages 4471–4480, 2018.
- Babak Esmaeili, Hao Wu, Sarthak Jain, N Siddharth, Brooks Paige, and Jan-Willem van de Meent. Hierarchical Disentangled Representations. *Artificial Intelligence and Statistics*, 2019.
- Niall P. Hurley and Scott T. Rickard. Comparing measures of sparsity. *IEEE Transactions on Information Theory*, 55:4723–4741, 2008.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning*, 2018.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Soheil Kolouri, Phillip E. Pope, Charles E. Martin, and Gustavo K. Rohde. Sliced wasserstein auto-encoders. In *International Conference on Learning Representations*, 2019.
- Chris J Maddison, John Lawson, George Tucker, Nicolas Heess, Mohammad Norouzi, Andriy Mnih, Arnaud Doucet, and Yee Teh. Filtering variational objectives. In *Advances in Neural Information Processing Systems*, pages 6573–6583, 2017.
- Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
- Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *International Conference on Learning Representations*, 2016.
- Tom Rainforth, Robert Cornish, Hongseok Yang, Andrew Warrington, and Frank Wood. On nesting monte carlo estimators. In *International Conference on Machine Learning*, pages 4264–4273, 2018a.
- Tom Rainforth, Adam R. Kosiorek, Tuan Anh Le, Chris J. Maddison, Maximilian Igl, Frank Wood, and Yee Whye Teh. Tighter variational bounds are not necessarily better. *International Conference on Machine Learning*, 2018b.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schölkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.