

Benefiting Deep Latent Variable Models via Learning the Prior and Removing Latent Regularization

Rogan Morrow and Wei-Chen Chiu

National Chiao Tung University, Taiwan
 rogan.o.morrow@gmail.com walon@cs.nctu.edu.tw

Abstract. There exist many forms of deep latent variable models, such as the variational autoencoder and adversarial autoencoder. Regardless of the specific class of model, there exists an implicit consensus that the latent distribution should be regularized towards the prior, even in the case where the prior distribution is learned. Upon investigating the effect of latent regularization on image generation our results indicate that in the case where a sufficiently expressive prior is learned, latent regularization is not necessary and may in fact be harmful insofar as image quality is concerned. We additionally investigate the benefit of learned priors on two common problems in computer vision: latent variable disentanglement, and diversity in image-to-image translation.

Keywords: VAEs, AAEs, generative autoencoders, disentanglement

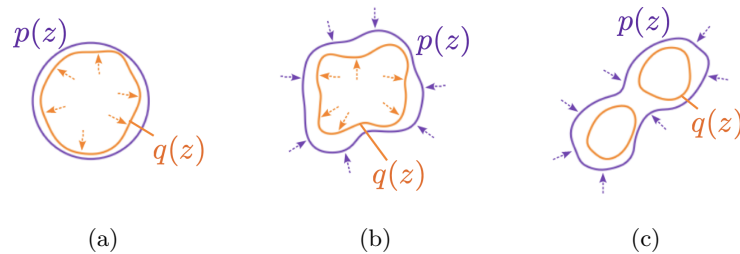


Fig. 1: Contour graph of prior distribution $p(\mathbf{z})$ and aggregated encoder distribution $q(\mathbf{z})$ for three different approaches to generative autoencoder training. Arrows represent forces acting on each distribution during training, excluding reconstruction loss. (a) Fixed $p(\mathbf{z})$, regularized $q(\mathbf{z})$. (b) Learned $p(\mathbf{z})$, regularized $q(\mathbf{z})$. (c) Learned $p(\mathbf{z})$, unregularized $q(\mathbf{z})$.

1 Introduction

In the machine learning subfield of deep latent variable models, generative autoencoders such as variational autoencoders (VAEs) [30] and adversarial autoencoders (AAEs) [38] have attracted a significant amount of research interest [4,29,20,45,7,51]. Despite this, in their standard form they are still largely outperformed in terms of synthesized image quality by other deep generative models such as generative adversarial networks (GANs) [12,25,26], autoregressive models [42] and flow-based models [8,9,31]. Even so, generative autoencoders maintain a number of properties that make them an attractive alternative, such as stable and efficient training as well as efficient synthesis.

In nearly all research done using such models, some form of regularization is imposed on the aggregated encoder distribution $q(\mathbf{z})$ in order to push it towards the prior distribution $p(\mathbf{z})$. For VAEs, this regularization exists in the form of a KL divergence between approximate posterior and prior, while AAEs force the aggregated posterior distribution to match the prior using an adversarial loss. This is of course necessary if the prior is fixed as is often the case, however we argue that when the prior is learnable it is possible to achieve a tight fit between aggregated posterior and prior without regularization, and that regularization in this case may actually have a negative impact on sample quality. Furthermore, removing regularization may result in a latent distribution that is beneficial to certain tasks such as disentanglement. Our contributions are as follows:

- We demonstrate empirically that when a sufficiently expressive prior $p(\mathbf{z})$ is learned, regularization of $q(\mathbf{z})$ is not necessary and may in fact be harmful to image quality.
- We demonstrate that when the linear disentanglement metrics proposed in [26] are considered, a learned prior outperforms other methods commonly used for generative autoencoder disentanglement, and that regularization of $q(\mathbf{z})$ does not improve linear disentanglement. This is in contrast to the common method of adding stronger regularization for $q(z)$ such as in [16].
- We demonstrate that a learned prior is beneficial to sample diversity in multimodal image-to-image translation tasks, where higher diversity is often a stated goal.

2 Generative Autoencoders and Latent Regularization

In this paper we focus on autoencoder-based generative models. This class of models defines an encoder distribution $q(\mathbf{z}|\mathbf{x})$ and a decoder distribution $p(\mathbf{x}|\mathbf{z})$, where the data $\mathbf{x} \sim p(\mathbf{x})$ is a random vector residing in space \mathcal{X} and \mathbf{z} is a latent code residing in space \mathcal{Z} . The negative log of $p(\mathbf{x}|\mathbf{z})$ is often referred to as the *reconstruction loss*. Let $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ be a set of i.i.d. observations drawn from the data distribution. Then the objective is to maximize

$$\sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}^{(i)})} [\log p(\mathbf{x}^{(i)}|\mathbf{z})] - R(\mathbf{x}^{(i)}) \quad (1)$$

where $R(\mathbf{x})$ is a regularization term. Defining a prior $p(\mathbf{z})$ allows us to generate samples from the model by first sampling $\hat{\mathbf{z}} \sim p(\mathbf{z})$ and then sampling $\hat{\mathbf{x}} \sim p(\mathbf{x}|\hat{\mathbf{z}})$. Clearly in order for the generative distribution of the model to closely match $p(\mathbf{x})$, the aggregated encoder distribution $q(\mathbf{z}) = \int_{\mathcal{X}} q(\mathbf{z}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$ should closely match the prior, such that we have $q(\mathbf{z}) \approx p(\mathbf{z})$ for all $\mathbf{z} \in \mathcal{Z}$. Therefore the regularization term $R(\mathbf{x})$ should be defined in such a way that it pushes $q(\mathbf{z})$ towards $p(\mathbf{z})$. Typically $p(\mathbf{z})$ is fixed as e.g. a standard normal distribution, however it is also possible to use the regularization term to learn the parameters of $p(\mathbf{z})$. If the prior is sufficiently expressive it may even be possible to remove regularization of $q(\mathbf{z})$ entirely so that the induced distribution of $q(\mathbf{z})$ is determined solely by pressure from the reconstruction loss, and the divergence between $q(\mathbf{z})$ and $p(\mathbf{z})$ is minimized solely by learning $p(\mathbf{z})$. These different approaches are visualized in Figure 1. Examples of the first approach with fixed $p(\mathbf{z})$ are ubiquitous in the literature [30,29,38,51], indeed it would be possible to fill an entire page with references to previous works utilizing this approach. The second approach, with both learned $p(\mathbf{z})$ and regularized $q(\mathbf{z})$, is less common but still abundant [6,52,24]. The third approach with unregularized $q(\mathbf{z})$ is exceedingly rare however; the only previous works we are aware of that adopt this approach are [33,36,2,54], and other than [54] they do not explicitly discuss the benefit of the approach¹. This motivates the question: does such an approach have any benefits over the first two, and should it be more commonly used? Intuitively regularization should impede $q(\mathbf{z})$ from assuming a shape that is most beneficial to the decoder, and so one would assume that reconstruction loss would be negatively affected. We discuss other potential model-specific downsides in the following subsection.

2.1 Related Models

Variational Autoencoders [30]

The goal of a likelihood-based model is to maximize the likelihood

$$p(\mathbf{x}^{(i)}) = \int_{\mathcal{Z}} p(\mathbf{x}^{(i)}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (2)$$

The integral in Eq. 2 is typically intractable, however. Variational autoencoders circumvent the issue of intractability by optimizing the variational lower bound

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}^{(i)})}[\log p(\mathbf{x}^{(i)}|\mathbf{z})] - D_{KL}[q(\mathbf{z}|\mathbf{x}^{(i)})||p(\mathbf{z})] \leq \log p(\mathbf{x}^{(i)}) \quad (3)$$

Thus we have $R(\mathbf{x}^{(i)}) = D_{KL}[q(\mathbf{z}|\mathbf{x}^{(i)})||p(\mathbf{z})]$. It is a well known issue that VAEs tend to not make full use of the latent code, as the objective becomes trapped in

¹ In [54], the authors focus specifically on VAE models and arrive at virtually the same result as us for our VAE derived model. However, their treatment of regularization and corresponding theoretical interpretation is fundamentally different to ours – they experiment on the value of the inverse variance of the decoder distribution, and interpret their unregularized model as the vanishing noise limit of a VAE.

a local minima in which the posterior is close to the prior, a phenomenon known as “posterior collapse” [13]. Such a state occurs early on when the signal from the latent code is weak, resulting in a weak reconstruction term that is easily outweighed by the KL divergence term. This causes the posterior distributions of data points to overlap such that the optimal decoding becomes a weighted mean of the data points in pixel space, typically resulting in blurry reconstructions. This issue is particularly pernicious in the conditional setting if care is not taken, as it is easily possible for the model to entirely ignore the latent code when it is conditioned on a relevant context, resulting in a deterministic mapping.

Many methods for encouraging use of the latent code have been proposed. For instance, annealing the KL divergence term from 0 to full strength [4,19] allows the model to largely ignore the KL divergence term at the beginning of training. “Free bits”, introduced in [29], places a limit on the information in nats per latent subset that can contribute to the KL divergence term, ensuring that each subset can contribute at least λ nats of information without penalty. In the context of conditional variational autoencoders, [58] proposed to add a latent reconstruction term to the objective to encourage the model to make full use of the latent code. All of these techniques are intended as a means of alleviating over-regularization imposed by the KL divergence term in the objective. If the prior is learned, however, it may be possible to eliminate such regularization entirely, hence obviating the need for any aforementioned techniques.

Adversarial Autoencoders [38]

Adversarial training [12] allows a distribution to be learned by playing a min-max game between a generator and a discriminator. The generator produces fake samples with the goal of fooling the discriminator, and the discriminator attempts to accurately classify samples as either real or fake. In the originally proposed setting where the discriminator outputs a probability, at optimality the model minimizes the Jensen-Shannon divergence between the data and generative distributions. Adversarial autoencoders apply this idea by using an adversarial term as the regularizer in Eq. 1. The discriminator is trained separately to maximize

$$\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log D(\mathbf{z})] + \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}^{(i)})}[\log(1 - D(\mathbf{z}))] \quad (4)$$

where D is the discriminator network. Note that $R(\mathbf{x}^{(i)})$ is equal to Eq. 4.

Regularization of the autoencoder in this way may introduce substantial noise. The reasoning for this is that discriminators are known to constantly shift their probability mass around during training in response to the generator, and so the decoder will be forced to deal with noisy latent codes. Removing the adversarial term from the autoencoder objective and instead using it to learn the prior dispels any such noise injection.

3 The Unregularized Generative Autoencoder Objective and its Connections to Optimal Transport

Removal of the regularization term in Eq. 1 presents an immediate problem – we desire to rely solely on the learning of $p(\mathbf{z})$ to minimize the divergence between $p(\mathbf{z})$ and $q(\mathbf{z})$, however $p(\mathbf{z})$ relies on the regularization term to learn. Therefore we must reformulate the objective so that learning of the prior is possible while the latent distribution remains unregularized. This can be achieved by casting the objective as a bilevel optimization problem. Let the encoder, decoder and prior distributions be denoted by $q_\phi(\mathbf{z}|\mathbf{x})$, $p_\psi(\mathbf{x}|\mathbf{z})$ and $p_\theta(\mathbf{z})$, and parameterized by ϕ , ψ and θ respectively. Then the objective becomes

$$\max_{\phi, \psi} F(\phi, \psi, \theta) \quad (5a)$$

$$\text{s.t. } \theta \in \arg \max_{\theta} f(\phi, \theta) \quad (5b)$$

where F and f are the upper and lower-level objectives and are given by

$$F(\phi, \psi, \theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\psi(\mathbf{x}^{(i)}|\mathbf{z})] - \beta R(\mathbf{x}^{(i)}; \phi, \theta) \quad (6)$$

$$f(\phi, \theta) = \sum_{i=1}^n -R(\mathbf{x}^{(i)}; \phi, \theta) \quad (7)$$

where we have introduced a hyperparameter β to control the strength of the regularization term; setting $\beta = 0$ allows us to remove regularization entirely without affecting learning of the prior. Note that when $\beta = 1$ the objective is equivalent to Eq. 1, therefore our approach is consistent with the original objective. The problem in Eq. 5 can be optimized straightforwardly via simultaneous gradient ascent by updating ϕ , ψ and θ using $\frac{\partial F}{\partial \phi}$, $\frac{\partial F}{\partial \psi}$ and $\frac{\partial f}{\partial \theta}$ respectively [37]. Consider the case where $\beta = 0$ and the summation over $R(\mathbf{x})$ in f corresponds to a divergence measure. If $p_\theta(\mathbf{z})$ is expressive enough to match any induced $q_\phi(\mathbf{z})$, then Eq. 5 is equivalent to the objective in the main theorem of [51], in which the authors demonstrate the equivalence between the optimal transport objective and

$$\min_{\phi, \psi} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [c(\mathbf{x}, G_\psi(\mathbf{z}))]] \quad (8a)$$

$$\text{s.t. } \mathcal{D}_{\mathbf{z}}(q_\phi(\mathbf{z})||p_\theta(\mathbf{z})) = 0 \quad (8b)$$

where $\mathcal{D}_{\mathbf{z}}$ is an arbitrary divergence measure, c is a cost function and G_ψ is a deterministic mapping $\mathcal{Z} \rightarrow \mathcal{X}$. The cost function is defined implicitly in Eq. 6 through the reconstruction loss, for example

$$p_\psi(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|G_\psi(\mathbf{z}), \mathbf{I}) \quad (9)$$

gives us the L2 cost function

$$c(\mathbf{x}, G_\psi(\mathbf{z})) = \|\mathbf{x} - G_\psi(\mathbf{z})\|_2^2 \quad (10)$$

plus a normalizing constant, in which case the model is performing 2-Wasserstein distance minimization between the data and generative distributions. When sampling from the generative distribution, in order to be consistent with the optimal transport interpretation it becomes necessary to use the deterministic mapping $G_\psi(\mathbf{z})$ rather than sampling from the full decoder distribution $p_\psi(\mathbf{x}|\mathbf{z})$, however this is already common practice when using decoders with simple distributions such as isotropic Gaussians. In practice, $p_\theta(\mathbf{z})$ will typically not be expressive enough to exactly match any induced $q_\phi(\mathbf{z})$, in which case Eq. 5 can be viewed as a relaxation of the constraint in Eq. 8. In [51] the authors also propose a relaxed objective by removing the constraint and adding a penalty to the objective, which coincides with Eq. 1. Note that when $R(\mathbf{x})$ is derived from the AAE objective, \mathcal{D}_z corresponds to the Jensen-Shannon divergence, whereas when $R(\mathbf{x})$ is derived from the VAE objective, \mathcal{D}_z corresponds to the Kullback-Leibler divergence.

4 Effect of Regularization: An Experimental Setup

Our goal here is to answer the question: does regularization of $q(\mathbf{z})$ hurt or help image quality when the prior is learned? We could simply compare two models, one with full regularization and one without any regularization, however it is possible that image quality as a function of regularization strength is not monotonic, and so we investigate how image quality varies with the strength of regularization by running experiments across varying values of β . We consider two kinds of models in our approach: VAEs with prior learned via normalizing flow, and AAEs with prior learned by a simple MLP. Normalizing flows [22] optimize a composition of bijective functions $f = f_T \circ \dots \circ f_1$ using the change of variables formula $p(\mathbf{u}) = p(\mathbf{h}) |\det(\frac{\partial \mathbf{h}}{\partial \mathbf{u}^T})|$, where $\mathbf{h} = f(\mathbf{u})$ and $p(\mathbf{h})$ is typically a standard normal distribution. Using a normalizing flow to learn the prior of a VAE was first proposed in [6], where the authors use an autoregressive flow. In our experiments we use a flow with affine coupling layers as proposed in [9]. This gives rise to a set of latents $\{\mathbf{z}_t \in \mathcal{Z}_t\}_{t=0}^T$, where $\mathbf{z}_T \sim q_\phi(\mathbf{z})$ and \mathbf{z}_0 should approximate $\mathcal{N}(\mathbf{0}, \mathbf{I})$ after training. In [6] the authors optimize a single objective

$$\sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\psi(\mathbf{x}^{(i)}|\mathbf{z}) + \beta(\log p_\theta(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}^{(i)}))] \quad (11)$$

Reformulating the objective so that it conforms to the bilevel structure of Eq. 5 we have

$$F(\phi, \psi, \theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\psi(\mathbf{x}^{(i)}|\mathbf{z}) + \beta(\log p_\theta(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}^{(i)}))] \quad (12)$$

$$f(\phi, \theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{z})] \quad (13)$$

For VAEs with a fixed prior (e.g. a standard normal distribution), as β becomes larger we achieve a more structured latent space at the expense of reconstruction

quality [16], and vice versa as β becomes smaller. When the parameters of the prior $p_\theta(\mathbf{z})$ are learnable, however, decreasing β does not necessarily sacrifice latent structure, as any additional incurred divergence between the aggregate posterior $q_\phi(\mathbf{z})$ and the prior $p_\theta(\mathbf{z})$ can be mitigated by adjusting the parameters of $p_\theta(\mathbf{z})$. In our experiments we consider VAEs with reconstruction loss given by an isotropic Gaussian, i.e. $p_\psi(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|G_\psi(\mathbf{z}), \gamma\mathbf{I})$ where $G_\psi(\mathbf{z})$ is the output of the decoder. It is common in many VAE implementations to keep the variance fixed, i.e. γ is fixed to a predetermined value and not altered during training. In [7] the authors propose learning the variance of the decoder distribution, and prove that it is always possible to achieve a better VAE cost by lowering the value of γ . We therefore consider both fixed $\gamma = 1$ and learned γ approaches in our experiments. Note that γ in effect, similarly to β , changes the strength of the regularization term. As γ becomes smaller the decoder distribution becomes more peaked, and so when the decoder has a good estimate of the mean the reconstruction loss will be much stronger relative to the KL divergence term. When γ is learned, the decoder is incentivized to lower the value of γ whenever it obtains a better estimate of the mean, and so we can expect that β will have less effect on the sample quality of the model than when γ is fixed. Learning the prior becomes necessary in this case, however, as the model will prioritize learning the data manifold over learning the ground truth distribution as pointed out in [7].

AAEs with learnable priors were first proposed in [24]. As we did for the VAE model, we again split the objective in order to conform with the bilevel objective in Eq. 5. Formally, we have the following objectives:

$$F(\phi, \psi, \theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\psi(\mathbf{x}^{(i)}|\mathbf{z}) - \beta \log(1 - D_\omega(\mathbf{z}))] \quad (14)$$

$$f(\phi, \theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z})} [-\log D_\omega(\mathbf{z})] \quad (15)$$

$$g(\omega, \phi, \theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z})} [\log D_\omega(\mathbf{z})] + \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log(1 - D_\omega(\mathbf{z}))] \quad (16)$$

where D_ω denotes the discriminator and ω denotes the parameters of the discriminator network. The discriminator objective is represented by g , and is optimized by $\arg \max_\omega g(\omega, \phi, \theta)$. These three objectives can be maximized iteratively. Sampling from the prior is performed by first sampling $u \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and then passing the sample through a simple MLP parameterized by θ .

5 Experiments on values of β

We first train on MNIST a VAE with normalizing flow prior and β set to zero with a 2-dimensional latent code in order to demonstrate visually the learned latent distributions. This is shown in Figure 2; while the autoencoder has learned

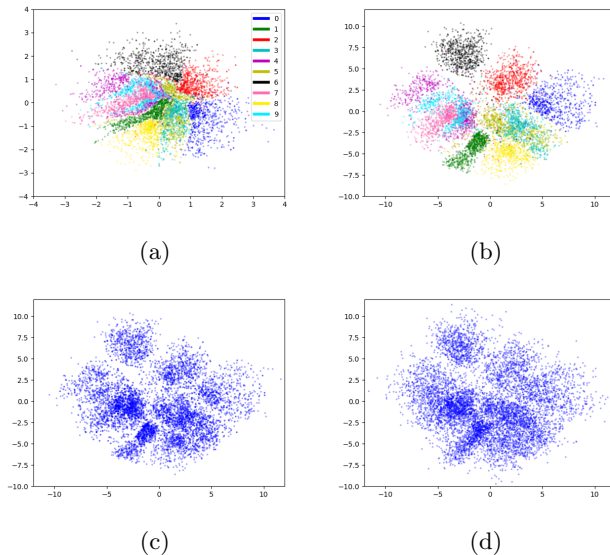


Fig. 2: Distribution of the latents of a VAE with flow prior with $\beta = 0$ after training on MNIST with $\dim(\mathcal{Z}) = 2$. (a) Distribution of $f(\mathbf{z})|_{\mathbf{z} \sim q(\mathbf{z})}$. Datapoints are colored according to class. (b) Distribution of $q(\mathbf{z})$. (c) Same as (b) but without class coloring. (d) Distribution of $f^{-1}(\mathbf{z})|_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})}$. It can be seen that while the autoencoder learns a complex latent distribution with classes well separated, the normalizing flow is able to learn a close match.

a latent distribution that is complex and multi-modal, the samples from the learned prior are a close match.

We experimented with varying values of β on the MNIST, Fashion-MNIST, CIFAR-10 and CelebA datasets. We chose to adopt the Fréchet Inception Distance (FID) [15] to measure image quality, a common measure used in GAN evaluation, for our quantitative comparisons. FID scores are given by the Fréchet distance between layer activations of the Inception v3 network [49], with lower scores indicating greater similarity between two image sets.

Results are reported in Figure 3. We report the average across 5 runs, and random seeds were kept fixed between runs such that the only changing hyperparameter is β . It can be seen that FID scores typically decrease as β is decreased for all models, suggesting that regularization of $q(z)$ is unnecessary and in fact potentially harmful to image quality. Interpolations in latent space for a VAE with L2 decoder and β set to zero are shown in Figure 4 in order to demonstrate that the model has learned a smooth manifold. We use spherical linear interpolation as suggested by [53].

FID scores were calculated against test sets using 10,000 samples. When calculating FID we used exactly the same code as was used in [7]. Although there may be slight discrepancies between different implementations of FID score, we

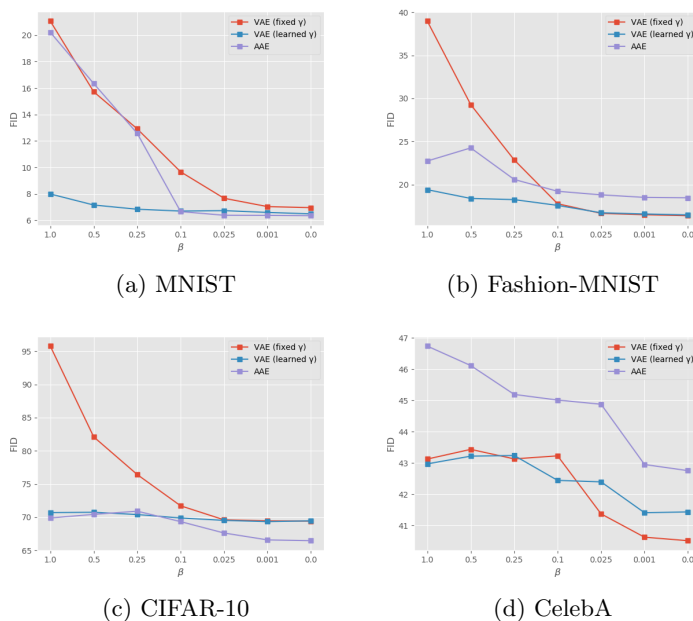


Fig. 3: Experimental results using different values of β for different models. X-axis represents β , not to scale. Y-axis represents FID score.

stress that our experiments are meant to be self-contained: we are primarily concerned with how β affects the same model trained under the same conditions, rather than how our results compare with those in other papers. We used a latent dimensionality of 64 for all datasets, and the same network architecture as was used in [5]. We note that we intentionally used priors that were of sufficient expressiveness to learn the latent distribution. For normalizing flows especially, this can result in a significant number of parameters being added to the model. We acknowledge that in cases where a lightweight model is desired or required such that a high capacity prior is not feasible, our results are not applicable, as latent regularization may still be required in order to achieve a good fit between encoder and prior distributions. For more details regarding the model architecture and training details, please see the supplementary.

Since there is no regularization imposed on the latent distribution at all when $\beta = 0$, it is possible that the dimensionality of the latent space becomes a critical hyperparameter when tuning the model. This is because it may be necessary to create an information bottleneck to induce a latent distribution that allows for the model to generalize well. An information bottleneck is also desirable for inducing a latent distribution that is easy enough for the prior network to learn. We experimented with how different values for the latent dimensions affects FID score for a VAE when $\beta = 0$, results are shown in Figure 5. We also included

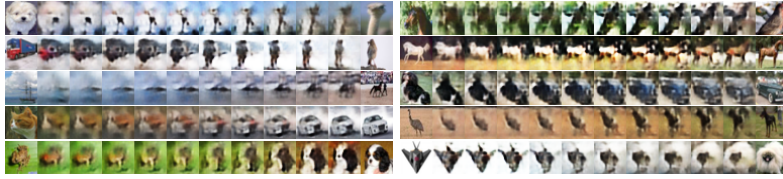


Fig. 4: Interpolation between samples from the CIFAR-10 test set for a flow prior with fixed γ and $\beta = 0$. Leftmost and rightmost columns contain real images from the test set before encoding, middle columns contain interpolations between them.

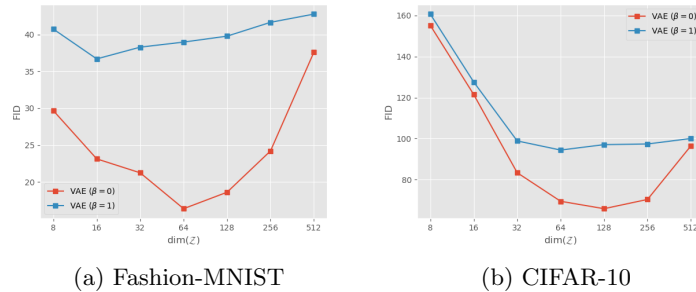


Fig. 5: Experimental results using different values for the latent dimensionality. X-axis represents the latent dimensionality. Y-axis represents FID score.

results for when $\beta = 1$ for reference. Note that the exact same prior network settings from the previous experiments was used, which was tuned for a latent dimensionality of 64. Results indicate that while lower values of β may achieve better image quality, it may be necessary to carefully tune the prior network and latent dimensionality, whereas standard VAEs tend to be more robust when considering relative change in FID score.

6 Disentanglement

The term “disentanglement” can cover a broad range of definitions, but a generalized high-level notion is that the model should capture individual factors of variation within linear subspaces of \mathcal{Z} . A common goal in the disentanglement literature is to regularize the model in such a way that it automatically aligns factors of variation along the axes of the latent space in an unsupervised manner [16,27]. Achieving this goal would therefore mean it is possible to perform semantic manipulation along an individual factor of variation by interpolating along any given axis in \mathcal{Z} . Clearly an unregularized autoencoder cannot achieve the same outcome, as it has no incentive to align factors of variation with the axes of \mathcal{Z} , and so extra processing is required in order to discover the directions

along which factors of variation lie. Despite this lack of automatic discovery, we argue that an unregularized $q(\mathbf{z})$ with a learned $p(\mathbf{z})$ has numerous benefits over existing disentanglements methods, which we discuss in the next two subsections.

6.1 Linear Disentanglement

Although an unregularized $q(\mathbf{z})$ does not allow automatic discovery of the factors of variation, we argue that it can achieve greater overall *linear disentanglement*. This means that, for any given semantic attribute, it should be easier to find a linear hyperplane that separates the latent codes into two sets, with each side of the hyperplane corresponding to one of the two possible values of the given attribute. Intuitively, a highly entangled latent representation cannot achieve good linear disentanglement, as it should not be possible to find a linear path in the latent space that can be interpolated along without altering many factors of variation simultaneously, and so semantic attributes cannot lie cleanly on either side of a linear hyperplane. As pointed out in [26], a fixed prior such as a standard normal necessarily entangles the latent space if there is any correlation between factors of variation. In [26] the authors posit that the decoder is likely to pressure $q(\mathbf{z})$ to take on a disentangled form, since intuitively this should make accurate reconstruction easier as opposed to trying to unwarped a highly entangled representation. Therefore it is reconstruction loss, not regularization, that induces better linear disentanglement.

We also argue that an unregularized $q(\mathbf{z})$ can achieve improved image quality *and* improved disentanglement simultaneously. Following from the point made above, reducing or removing regularization can only be beneficial to disentanglement, as it is the reconstruction loss that induces a disentangled representation. And, given that we have shown in Section 5 that removing regularization is beneficial to image quality, the two outcomes can be achieved simultaneously.

The CelebA dataset contains 40 binary attributes that we can consider as factors of variation, and thus we can use these attributes to calculate a measure of disentanglement. We consider the linear separability score proposed in [26]. In their work they first train a deep network classifier that predicts image attributes on the training images, and then train a linear SVM classifier that predicts the classifier network’s output given the latent variable. After this they calculate the conditional entropy $H(\mathbf{Y}|\mathbf{X})$ where \mathbf{Y} represents the labels predicted by the deep network classifier, and \mathbf{X} represents the labels predicted by the SVM. It can be seen that lower conditional entropy will correspond to better linear separation, since the SVM will have higher prediction accuracy and thus observing \mathbf{Y} will give less information. By following their procedure exactly, we can quantitatively measure linear disentanglement purely as a function of the generative process of the model. We additionally consider the perceptual path length proposed in [26]. This gives us a further measure of disentanglement; if the factors of variation lie along paths that are highly warped and curved, then a small movement along a linear (or spherical) path between two randomly sampled endpoints is likely to cause a larger perceptual change than if the factors of variation were lying along

linear paths.

We evaluated disentanglement in a VAE with standard normal prior with $\beta = 1$, a β -VAE [16] with $\beta = 25$, a FactorVAE [27] with the total correlation strength set to 40, a VAE with normalizing flow prior with $\beta = 1$ and a VAE with normalizing flow prior with $\beta = 0$. For both VAEs with normalizing flow prior we evaluated disentanglement of both \mathbf{z}_0 and \mathbf{z}_T to demonstrate that the fixed base distribution causes the latents to become significantly warped and entangled. Results are reported in Table 1. The results indicate that when $q(\mathbf{z})$ is unregularized and not fit to a fixed prior, linear disentanglement is improved, and additionally sample image quality improves. We provide samples from each model in Figure 7 as well as an example of feature discovery by performing PCA on $q(\mathbf{z})$. Clearly PCA is a poor method for feature discovery, however we believe it suffices for this simple demonstration.



Fig. 6: Adding glasses to a face. Top row of each set: interpolation in \mathcal{Z}_0 . Bottom row of each set: interpolation in \mathcal{Z}_T . Leftmost column: the original image before encoding. The top row shows a more abrupt change towards the end, while the bottom row is closer to a constant rate of change.

6.2 Interpolation

When interpolating between points in latent space the motivation is often to achieve some semantic mixture between two images, or to change some semantic feature of an image such as putting glasses on a person’s face. As discussed in the previous section, poor linear disentanglement is clearly detrimental to this task, as interpolation along a warped latent representation makes it difficult to manipulate a single semantic feature without potentially altering several other unrelated features. Issues of entanglement aside, we expect an unregularized $q(\mathbf{z})$ will achieve a more constant rate of change when interpolating along the direction of a factor of variation, especially in the case where the factor of variation

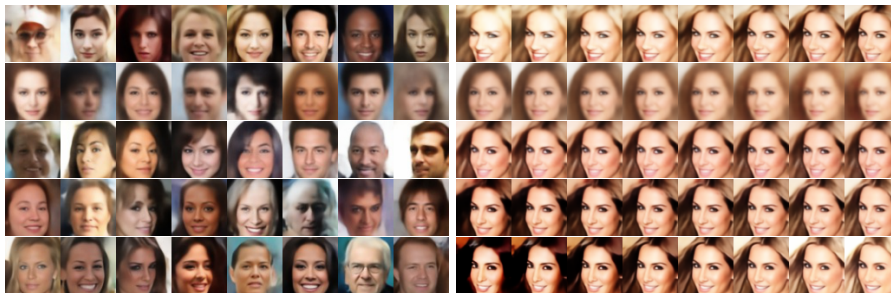


Fig. 7: Left columns: random samples from each model. Right columns: PCA was performed on $q(z)$ and an image from the dataset was encoded and transformed to PC space. The first principal component was interpolated along $[-1.5, 1.5]$ while other components were kept fixed. Row 1: VAE (std. normal prior). Row 2: β -VAE. Row 3: FactorVAE. Row 4: VAE (flow prior, $\beta = 1$). Row 5: VAE (flow prior, $\beta = 0$).

corresponds to a semantic attribute with class imbalance. This is because when using a fixed prior the rate of change in density is fixed according to the chosen prior distribution, and so the measure of variation is highly unlikely to have a linear correlation with distance travelled along the path. As an example, consider the case of a uniform prior distribution: the relative amount of space occupied by points with a particular semantic feature would be proportional to the class probability of the feature. What this means in practice is that if we were to generate a sequence of images by interpolating along the direction of a semantic attribute with heavy class imbalance, there would be very little change for most of the sequence followed by an abrupt change at the end.

A normalizing flow prior allows us to perform a fair comparison between interpolation in an unregularized distribution and interpolation in a fixed distribution, since the the bijection allows us to interpolate between corresponding path endpoints in either distribution. We therefore experiment on the difference between interpolating in \mathcal{Z}_0 or in \mathcal{Z}_T using a VAE with normalizing flow prior and $\beta = 0$. As a point of clarification, when we say we are interpolating between two images $\mathbf{x}^{(a)}$ and $\mathbf{x}^{(b)}$ in \mathcal{Z}_T we are calculating $\hat{\mathbf{x}} \sim p(\mathbf{x}|\hat{\mathbf{z}})$ where $\hat{\mathbf{z}} = \text{lerp}(\mathbf{z}^{(a)}, \mathbf{z}^{(b)}; t)$, $\mathbf{z}^{(a)} \sim q(\mathbf{z}|\mathbf{x}^{(a)})$, $\mathbf{z}^{(b)} \sim q(\mathbf{z}|\mathbf{x}^{(b)})$ and t varies between 0 and 1. When we are interpolating in \mathcal{Z}_0 we are instead calculating $\hat{\mathbf{z}} = f^{-1}(\text{slerp}(f(\mathbf{z}^{(a)}), f(\mathbf{z}^{(b)}); t))$, where f is the bijection defined by the normalizing flow. In order to calculate the direction of change for a particular attribute, we first calculate the mean of all latent codes \mathbf{z}_T corresponding to images with and without the attribute. We then calculate the difference between these two means to produce the direction of change. If the attribute corresponds to glasses, for example, we can encode an image of a person not wearing glasses, add the vector representing the direction of change to the latent encoding, and then decode to produce an image of the same person wearing glasses. We used VGG19 perceptual loss for the decoder of the model as we found this helped with

semantic manipulation using vector arithmetic. To quantitatively measure the rate of change, we sample 5000 images from the data set without the attribute, and then interpolate evenly along the direction of change to produce a sequence of 16 images. We then measure perceptual difference between adjacent images using a VGG16 network. Results are shown in Figure 8, where we plot the median perceptual change along the generated sequences. It can be seen that the rate of change when interpolating in \mathcal{Z}_0 varies significantly, especially at the ends of the sequence. Interpolation in \mathcal{Z}_T on the other hand produces a rate of change that is much closer to constant, which is more ideal for semantic manipulation. Some sample sequences are shown in Figure 6.

For architecture and training details regarding our disentanglement experiments, please refer to the supplementary.

Table 1: Linear separability, perceptual path length (PPL) and FID scores after training on CelebA.

	Separability	PPL	FID
VAE (std. normal prior)	2.14	1139	41.4
β -VAE	2.09	1324	82.7
FactorVAE	2.32	1339	39.8
VAE (flow prior, $\beta = 1$) (\mathbf{z}_0)	2.46	1753	38.7
VAE (flow prior, $\beta = 1$) (\mathbf{z}_T)	1.68	1173	
VAE (flow prior, $\beta = 0$) (\mathbf{z}_0)	2.82	1933	33.1
VAE (flow prior, $\beta = 0$) (\mathbf{z}_T)	1.67	1076	

Table 2: Inception and LPIPS scores (higher is better) on the DeepFashion dataset after training Variational U-net.

	IS	LPIPS
VUNET (Original)	2.63	0.184
VUNET (Proposed)	2.70	0.236

7 Diversity in image-to-image translation

Achieving high sample diversity in multi-modal image-to-image translation tasks is often an explicit goal [58,21]. When using conditional VAEs for image-to-image translation tasks, the decoder is often able to learn a fairly accurate reconstruction based on the conditioned image alone, and so may ignore the latent code entirely if the KL divergence weight is too strong. In order to quantitatively

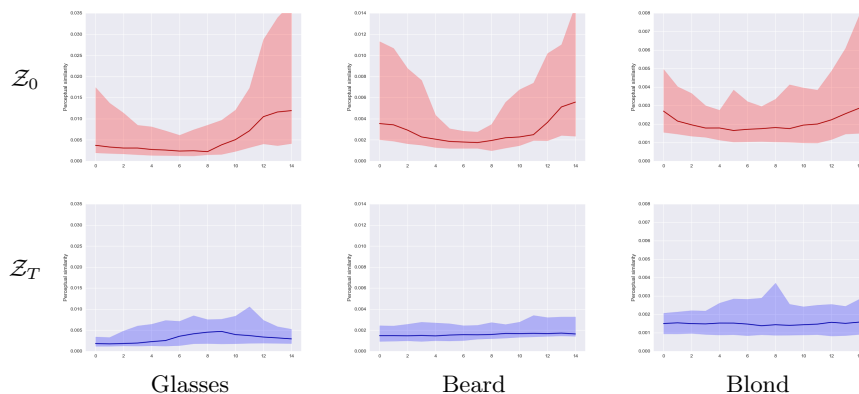


Fig. 8: We plot the median perceptual change for each point in multiple sequences generated by manipulating a particular semantic feature in either \mathcal{Z}_0 or \mathcal{Z}_T . The shaded region represents the 5th to 95th percentiles. X-axis represents index in the sequence. Y-axis represents perceptual change.

test whether our proposed method is able improve diversity, we experiment with the Variational U-net model proposed in [10]. In their work, they attempt to learn the distribution over images of people conditioned on their pose. In their implementation they make use of KL divergence annealing in order to encourage the model to make use of the latents, however sample diversity may still be negatively affected. We modified their implementation such that the prior distribution is learned via normalizing flow, and dropped the KL divergence term from the objective. Their model conditions the prior distribution on the given pose such that their objective becomes

$$\log p(\mathbf{x}^{(i)}|\mathbf{y}^{(i)}, \mathbf{z}) - D_{KL}[q(\mathbf{z}|\mathbf{x}^{(i)}, \mathbf{y}^{(i)})||p(\mathbf{z}|\mathbf{y}^{(i)})] \quad (17)$$

where \mathbf{y} is the pose and \mathbf{x} is the real image. For compatibility we therefore learn the mean of \mathbf{z}_0 conditioned on the pose and additionally condition each flow transformation f_t on the pose. To measure diversity, we compute the LPIPS distance [56] between randomly sampled pairs which were generated by conditioning on the images in the test set. We additionally calculated the Inception score [46] of the samples to ensure image quality was not affected. Results comparing our modification with the original implementation after training on the DeepFashion dataset are reported in Table 2. We also show samples in Figure 9.

8 Conclusion

We have proposed removing latent regularization from the objective of generative autoencoder models in the case where the prior is sufficiently expressive. We demonstrated empirically that this results in improved image quality, improved

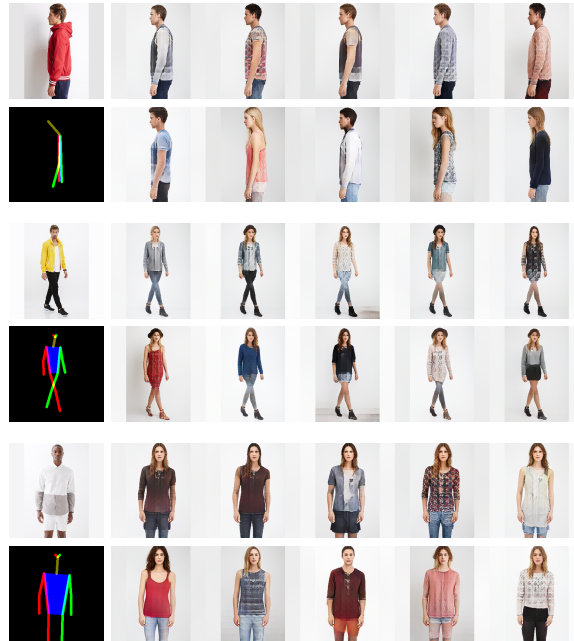


Fig. 9: Conditional samples using Variational U-net. For each set, the top row contains samples from the original, while the bottom row contains samples from the proposed change. Leftmost column contains the original image and the pose being conditioned on.

linear disentanglement, and improved sample diversity. Our results indicate that fixed-form priors should be eschewed in favour of learned priors with little to no latent regularization.

References

1. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein GAN. In: International Conference on Machine Learning (ICML) (2017)
2. Bojanowski, P., Joulin, A., Lopez-Pas, D., Szlam, A.: Optimizing the latent space of generative networks. In: Proceedings of the 35th International Conference on Machine Learning. pp. 600–609 (2018)
3. Borji, A.: Pros and cons of GAN evaluation measures. *Computer Vision and Image Understanding (CVIU)* (2019)
4. Bowman, S.R., Vilnis, L., Vinyals, O., Dai, A.M., Józefowicz, R., Bengio, S.: Generating sentences from a continuous space. In: The SIGNLL Conference on Computational Natural Language Learning (CoNLL) (2016)
5. Chen, X., Duan, Y., Houthoofd, R., Schulman, J., Sutskever, I., Abbeel, P.: InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In: Advances in Neural Information Processing Systems (NIPS) (2016)

6. Chen, X., Kingma, D.P., Salimans, T., Duan, Y., Dhariwal, P., Schulman, J., Sutskever, I., Abbeel, P.: Variational lossy autoencoder. In: International Conference on Learning Representations (ICLR) (2017)
7. Dai, B., Wipf, D.: Diagnosing and enhancing VAE models. In: International Conference on Learning Representations (ICLR) (2019)
8. Dinh, L., Krueger, D., Bengio, Y.: NICE: non-linear independent components estimation. In: International Conference on Learning Representations (ICLR) (2015)
9. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real NVP. In: International Conference on Learning Representations (ICLR) (2017)
10. Esser, P., Sutter, E., Ommer, B.: A variational U-net for conditional appearance and shape generation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
11. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: Teh, Y.W., Titterton, M. (eds.) Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research, vol. 9, pp. 249–256. PMLR, Chia Laguna Resort, Sardinia, Italy (13–15 May 2010)
12. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in Neural Information Processing Systems (NIPS) (2014)
13. He, J., Spokoyny, D., Neubig, G., Berg-Kirkpatrick, T.: Lagging inference networks and posterior collapse in variational autoencoders. In: Proceedings of ICLR (2019)
14. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015)
15. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In: Advances in Neural Information Processing Systems (NIPS) (2017)
16. Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., Lerchner, A.: beta-VAE: Learning basic visual concepts with a constrained variational framework. In: International Conference on Learning Representations (ICLR) (2017)
17. Hou, X., Shen, L., Sun, K., Qiu, G.: Deep feature consistent variational autoencoder. In: IEEE Winter Conference on Applications of Computer Vision (WACV) (2017)
18. Huang, C., Krueger, D., Lacoste, A., Courville, A.C.: Neural autoregressive flows. In: International Conference on Machine Learning (ICML) (2018)
19. Huang, C., Tan, S., Lacoste, A., Courville, A.C.: Improving explorability in variational inference with annealed variational objectives. In: Advances in Neural Information Processing Systems (NIPS) (2018)
20. Huang, C., Touati, A., Dinh, L., Drozdal, M., Havaei, M., Charlin, L., Courville, A.C.: Learnable explicit density for continuous latent space and variational inference. In: International Conference on Machine Learning (ICML) Workshops (2017)
21. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: European Conference on Computer Vision (ECCV) (2018)
22. Jimenez Rezende, D., Mohamed, S.: Variational Inference with Normalizing Flows. In: International Conference on Machine Learning (ICML) (2015)
23. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: European Conference on Computer Vision (ECCV) (2016)
24. Junbo, Zhao, Kim, Y., Zhang, K., Rush, A.M., LeCun, Y.: Adversarially Regularized Autoencoders for Generating Discrete Structures. ArXiv e-prints (2017)

25. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. In: International Conference on Learning Representations (ICLR) (2018)
26. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. ArXiv:1812.04948 (2018)
27. Kim, H., Mnih, A.: Disentangling by Factorising. In: Advances in Neural Information Processing Systems (NIPS) Workshops (2017)
28. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization. arXiv e-prints (Dec 2014)
29. Kingma, D.P., Salimans, T., Welling, M.: Improving variational inference with inverse autoregressive flow. In: Advances in Neural Information Processing Systems (NIPS) (2016)
30. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: International Conference on Learning Representations (ICLR) (2014)
31. Kingma, D.P., Dhariwal, P.: Glow: Generative flow with invertible 1x1 convolutions. In: Advances in Neural Information Processing Systems (NIPS) (2018)
32. Li, C., Chang, W., Cheng, Y., Yang, Y., Póczos, B.: MMD GAN: towards deeper understanding of moment matching network. In: Advances in Neural Information Processing Systems (NIPS) (2017)
33. Li, Y., Swersky, K., Zemel, R.S.: Generative moment matching networks. In: International Conference on Machine Learning (ICML) (2015)
34. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: IEEE International Conference on Computer Vision (ICCV) (2015)
35. Lucic, M., Kurach, K., Michalski, M., Gelly, S., Bousquet, O.: Are GANs Created Equal? A Large-Scale Study. In: Advances in Neural Information Processing Systems (NIPS) (2018)
36. Ma, L., Sun, Q., Georgoulis, S., Van Gool, L., Schiele, B., Fritz, M.: Disentangled person image generation. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)
37. MacKay, M., Vicol, P., Lorraine, J., Duvenaud, D., Grosse, R.: Self-tuning networks: Bilevel optimization of hyperparameters using structured best-response functions. In: International Conference on Learning Representations (ICLR) (2019)
38. Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I.J.: Adversarial autoencoders. In: International Conference on Learning Representations (ICLR) (2016)
39. Mescheder, L.M.: On the convergence properties of GAN training. ArXiv:1801.04406 (2018)
40. Metz, L., Poole, B., Pfau, D., Sohl-Dickstein, J.: Unrolled generative adversarial networks. In: International Conference on Learning Representations (ICLR) (2017)
41. Odena, A., Olah, C., Shlens, J.: Conditional Image Synthesis With Auxiliary Classifier GANs. In: International Conference on Machine Learning (ICML) (2016)
42. Van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al.: Conditional image generation with pixelcnn decoders. In: Advances in Neural Information Processing Systems (NIPS) (2016)
43. Papamakarios, G., Pavlakou, T., Murray, I.: Masked Autoregressive Flow for Density Estimation. In: Advances in Neural Information Processing Systems (NIPS) (2017)
44. Razavi, A., van den Oord, A., Vinyals, O.: Generating diverse high-fidelity images with vq-vae-2. In: Advances in Neural Information Processing Systems 32. pp. 14866–14876 (2019)
45. Rezende, D.J., Viola, F.: Taming VAEs. ArXiv:1810.00597 (2018)

46. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training GANs. In: *Advances in Neural Information Processing Systems (NIPS)* (2016)
47. Salimans, T., Karpathy, A., Chen, X., Kingma, D.P.: Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. In: *International Conference on Learning Representations (ICLR)* (2017)
48. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: *International Conference on Learning Representations (ICLR)* (2015)
49. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015)
50. Theis, L., van den Oord, A., Bethge, M.: A note on the evaluation of generative models. *ArXiv:1511.01844* (2015)
51. Tolstikhin, I., Bousquet, O., Gelly, S., Schoelkopf, B.: Wasserstein Auto-Encoders. In: *International Conference on Learning Representations (ICLR)* (2018)
52. Tomczak, J.M., Welling, M.: VAE with a VampPrior. *arXiv* (2017)
53. White, T.: Sampling generative networks: Notes on a few effective techniques. *ArXiv:1609.04468* (2016)
54. Xiao, Z., Yan, Q., Chen, Y., Amit, Y.: Generative latent flow: A framework for non-adversarial image generation. *CoRR* **abs/1905.10485** (2019)
55. Xu, Q., Huang, G., Yuan, Y., Guo, C., Sun, Y., Wu, F., Weinberger, K.Q.: An empirical study on evaluation metrics of generative adversarial networks. *ArXiv:1806.07755* (2018)
56. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
57. Zhao, S., Song, J., Ermon, S.: Towards deeper understanding of variational autoencoding models. *ArXiv:1702.08658* (2017)
58. Zhu, J., Zhang, R., Pathak, D., Darrell, T., Efros, A.A., Wang, O., Shechtman, E.: Toward multimodal image-to-image translation. In: *Advances in Neural Information Processing Systems (NIPS)* (2017)

A Regularization experiment

We use the same autoencoder architecture used in [5] for all datasets.

A.1 Flow prior architecture

We use the same normalizing flow architecture to model the prior in each of our experiments. Here we give a detailed overview of the architecture. The type of flow we used is the same as that of RealNVP [9]. That is, for each individual transformation in the flow we split the input into two parts, and pass one part through a neural network to give the parameters of an affine transformation that is applied to the second part. We additionally apply an affine transformation with learnable parameters before every second transformation, the same as the `actnorm` operation used in [31]. We do not use data-dependent initialization for the affine transformation, and instead initialize it as the identity function.

z	\leftarrow	$\gamma z + \beta$	\triangleright affine transformation
z_a, z_b	\leftarrow	$\text{split}(z)$	
μ, σ	\leftarrow	$\text{NN}(z_a)$	
σ	\leftarrow	$\text{sigmoid}(\sigma + 2)$	
z_b	\leftarrow	$\sigma(z_b + \mu)$	
z	\leftarrow	$\text{concat}(z_a, z_b)$	

Pseudo-code for a single transformation f_t within the flow is given below:

Next we discuss the structure of the network NN. We use the following denotations to describe the architecture: FC-X denotes a fully connected layer with an X-dimensional output. Let D denote the “width” of a flow transform, and let R equal $\dim(\mathcal{Z})/2$. Then NN is given by FC-D \rightarrow ReLU \rightarrow FC-D \rightarrow ReLU \rightarrow FC-R.

All kernel weights in NN are initialized using Glorot uniform initialization [11], except for the final fully connected layer whose weights are initialized as zero.

A.2 Adversarial prior architecture

The adversarial prior architecture used across all experiments is as follows. Gaussian noise is passed through several layers of width D , followed by a final fully connected layer of width R . Let BN denote batch normalization. Each layer is given by FC-D \rightarrow BN \rightarrow ReLU.

The discriminator is also given by several layers of width D followed by a final fully connected layer of width 1. Each layer is given by FC-D \rightarrow ReLU.

A.3 Hyperparameters and training

A batch size of 100 and an initial learning rate of 0.0001 was used across all datasets and models. We used the Adam optimizer [28] with default parameters.

The details for the flow prior experiments are as follows.
 For MNIST we used a 24 layer flow with a width of 1024. The entire model was trained for 200 epochs, and the prior was then trained independently for a further 100 epochs.
 For Fashion-MNIST we used a 24 layer flow with a width of 1024. The entire model was trained for 300 epochs, and the prior was then trained independently for a further 100 epochs.
 For CIFAR-10 we used a 16 layer flow with a width of 256. The entire model was trained for 500 epochs, and the prior was then trained independently for a further 100 epochs. The learning rate was halved every 250 epochs.

For CelebA we used a 30 layer flow with a width of 1024. The entire model was trained for 200 epochs, and the prior was then trained independently for a further 50 epochs.

The details for the adversarial prior experiments are as follows. We used the same size prior and discriminator across all datasets; the prior was 3 layers of width 1024 followed by a final fully connected layer and the discriminator was 2 layers of width 1024 followed by a final fully connected layer.

For MNIST we trained the entire model for 200 epochs.

For Fashion-MNIST we trained the entire model for 300 epochs.

For CIFAR-10 we trained the entire model for 300 epochs.

For CelebA we trained the entire model for 200 epochs.

B Disentanglement experiments

For our disentanglement experiments using CelebA we used exactly the same Resnet architecture as described in [7] with a depth of 4. The prior architecture used for the VAEs with normalizing flow prior was the same as described in the previous section. L2 loss was used for the decoder. We trained for 300 epochs and halved the learning rate every 100 epochs.

B.1 VGG19 for interpolation

We used the squared difference between the hidden features of the `relu_1_1`, `relu_2_1`, `relu_3_1` and `relu_4_1` layers of the VGG19 network as the reconstruction loss for the interpolation experiment.