
On the Regularization Balance in Autoencoder-Based Generative Models

Anonymous Authors¹

Abstract

Although GANs have recently displayed impressive visual results, there remains value in generative models that do not rely on adversarial training and are capable of producing posterior estimates of latent representations when conditioned on an input. Within this scope, we analyze generative models of continuous data based upon a regularized autoencoder structure, e.g., variational autoencoders (VAEs) and related architectures. First, we demonstrate that with an affine decoder there will exist a combinatorial number of distinct local minima; however, each of these will necessarily be global optima displaying balanced regularization. By this we mean that perfect reconstructions can be obtained using a minimal number of latent factors, which we contend is essential for generating realistic samples with an autoencoder structure. We then contrast the balanced case with more complex situations where over-regularized solutions may produce poor reconstructions, or under-regularized solutions may lead to unnecessarily redundant latent representations. Proceeding further, we discuss practical solutions for addressing over-regularization, while arguing that under-regularization avoidance in a certain sense requires the use of energy functions with infinite gradients (or nearly so) around optima (e.g., a Gaussian VAE model with an optimal decoder variance near zero). While potentially leading to numerical issues, such energy functions are essentially unavoidable across a wide range of traditional optimization problems designed to produce maximum parsimony, but strategies exist to at least partially mitigate instabilities. Empirical examples using real-world images corroborate our findings and demonstrate that non-adversarial training, when properly regularized, can produce reasonable samples even with full-resolution CelebA data.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

1. Introduction

Suppose we have access to continuous variables \mathbf{x} that are assumed to lie on or near an r -dimensional manifold \mathcal{X} embedded in \mathbb{R}^d , where $r \ll d$. This construction captures the notion that our data of interest possesses low-dimensional structure relative to the high-dimensional ambient space. We denote a ground-truth probability measure on \mathcal{X} as μ_{gt} , which assigns probability mass $\mu_{gt}(d\mathbf{x})$ to the infinitesimal $d\mathbf{x}$ residing on the manifold; it naturally follows that $\int_{\mathcal{X}} \mu_{gt}(d\mathbf{x}) = 1$. Given samples $\{\mathbf{x}^{(i)}\}_{i=1}^n$ drawn from this measure, our goal is to estimate the unknown μ_{gt} , or at least to provide a tractable way of drawing new samples that are indistinguishable in distribution from the training sample.

Because of the assumed low-dimensional structure, it is common to approximate the unknown ground-truth measure via a model parameterized as $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$, where θ are parameters, $\mathbf{z} \in \mathbb{R}^{\kappa}$ serves as a low-dimensional representation, with $\kappa \approx r$ and fixed prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$. Samples from $p_{\theta}(\mathbf{x})$ can then be easily produced by drawing $\mathbf{z} \sim p(\mathbf{z})$ and then $\mathbf{x} \sim p_{\theta}(\mathbf{x}|\mathbf{z})$, noting that if $p_{\theta}(\mathbf{x}|\mathbf{z})$ happens to be a Dirac delta function, the only source of randomness arises from sampling \mathbf{z} , which is then processed by what defaults to a deterministic transformation to produce \mathbf{x} .

Recently, generative adversarial networks (GANs) have demonstrated impressive results generating novel signals/images by learning such a deterministic $p_{\theta}(\mathbf{x}|\mathbf{z}) \equiv \delta[\mathbf{x} - f_{\theta}(\mathbf{z})]$ (Goodfellow et al., 2014). The potential downside though is that the required adversarial training can be sensitive to mode collapse (Arora & Zhang, 2017; Metz et al., 2016), hyperparameter tuning (Lucic et al., 2018), and other implementation/training issues. Moreover, standard GAN models do not provide a tractable way of computing the posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$, which is useful for a variety of applications (Dai et al., 2018). However, there also exist modeling approaches that do not rely on adversarial training while producing direct posterior estimates. The variational autoencoder (VAE) and its many variants represent influential examples that fall within this category (Kingma & Welling, 2014; Rezende et al., 2014).

In principle we might consider minimizing the negative log-likelihood $-\log p_{\theta}(\mathbf{x})$ averaged across the training data, approximating integration across the assumed μ_{gt} , i.e., min-

imize

$$\frac{1}{n} \sum_i -\log [p_\theta(\mathbf{x}^{(i)})] \approx \int -\log [p_\theta(\mathbf{x})] \mu_{gt}(d\mathbf{x}) \quad (1)$$

over θ . Unfortunately though, the marginalization required to produce $p_\theta(\mathbf{x}^{(i)})$ is generally intractable for models of sufficient representational power. To circumvent this issue, the VAE instead attempts to minimize the variational upper bound $\mathcal{L}(\theta, \phi) \triangleq$

$$\begin{aligned} & \int \{-\log p_\theta(\mathbf{x}) + \mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})]\} \mu_{gt}(d\mathbf{x}) \\ &= \int \{-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \\ & \quad + \mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]\} \mu_{gt}(d\mathbf{x}) \quad (2) \end{aligned}$$

or a sampling-based approximation. Here $q_\phi(\mathbf{z}|\mathbf{x})$ represents a tractable variational approximation to $p_\theta(\mathbf{z}|\mathbf{x})$ with additional parameters ϕ governing the tightness of the bound. It is commonly referred to as an *encoder* distribution since it quantifies the mapping from \mathbf{x} to the latent code \mathbf{z} . For analogous reasons, $p_\theta(\mathbf{x}|\mathbf{z})$ is labeled as the *decoder* distribution. When combined, the data-dependent factor $-\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$ can be viewed as instantiating a form of stochastic autoencoder (AE) structure, which attempts to assign high probability to accurate reconstructions of each \mathbf{x} ; if $q_\phi(\mathbf{z}|\mathbf{x})$ is Dirac delta function, then a regular deterministic AE emerges with loss dictated by the decoder negative log-likelihood $-\log p_\theta(\mathbf{x}|\mathbf{z})$. Beyond this, $\mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ serves as a regularization factor that pushes the encoder distribution towards the prior. The bound $\mathcal{L}(\theta, \phi)$, or at least its approximate version given samples from μ_{gt} , can be minimized over $\{\theta, \phi\}$ using SGD and a simple reparameterization trick (Kingma & Welling, 2014; Rezende et al., 2014).

The VAE as described above is representative of a wider class of generative models built upon an AE structure. Broadly speaking, if the overriding objective is generating realistic samples using an encoder-decoder-based architecture (VAE or otherwise), two important, well-known criteria must be satisfied:

- (i) Small reconstruction error when passing through the encoder-decoder networks, and
- (ii) An *aggregate posterior* distribution (Makhzani et al., 2016) over latent codes \mathbf{z} , defined by $q_\phi(\mathbf{z}) \triangleq \int q_\phi(\mathbf{z}|\mathbf{x}) \mu_{gt}(d\mathbf{x})$, that is close to a known distribution like $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|0, I)$ that is easy to sample from.

The first criteria can be naturally enforced by a deterministic AE, but also for the VAE provided that sufficient randomness can be suppressed by the encoder. Of course the second criteria is equally important. Without it, we have no tractable way of generating random decoder inputs that, when passed to the decoder output, are transformed to realistic samples resembling the true data distribution. Additionally, if we assume that $\kappa = r$, then criteria (ii) involves comparing

distributions that should have finite density across a smaller r -dimensional space. This avoids potentially problematic comparisons in \mathbb{R}^d between measures that assign all or most probability mass to a lower-dimensional manifold with infinite density.

Criteria (i) can in principle be handled by a sufficiently complex aggregate AE network, while (ii) can be addressed by increasing the flexibility of $q_\phi(\mathbf{z}|\mathbf{x})$ (Burda et al., 2015; Kingma et al., 2016; Rezende & Mohamed, 2015; van den Berg et al., 2018) or the prior $p(\mathbf{z})$ (Bauer & Mnih, 2018; Tomczak & Welling, 2018) such that there are additional pathways for pushing the intractable aggregate posterior towards a tractable prior that we can sample from. Other approaches exist as well. For example, unlike the VAE which relies on a KL regularization factor to indirectly favor $q_\phi(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z})$, the Wasserstein autoencoder (WAE) (Tolstikhin et al., 2018) instead utilizes an energy function that more directly penalizes the distance between $q_\phi(\mathbf{z})$ and $p(\mathbf{z})$. The WAE can operate using either adversarial training akin to the adversarial autoencoder (Makhzani et al., 2016), or via a non-adversarial, maximum mean discrepancy (MMD) penalty as considered herein. In contrast, it has also been shown that a second-stage VAE module can be used to sample from $q_\phi(\mathbf{z})$, even when its distribution is far from $p(\mathbf{z})$ (Dai & Wipf, 2018).

Despite significant differences, all of these methods are sensitive to the regularization balance required for simultaneously addressing criteria (i) and (ii). Although alternative definitions are certainly possible, for present purposes we adopt a very specific notion of what constitutes *balanced regularization*. In particular, we will make the case that a necessary condition for achieving high-quality samples via an autoencoder structure devoid of adversarial training, is an encoder-decoder pair that produces high-quality reconstructions of the training data per criteria (i), using a minimal number of latent degrees-of-freedom, a prerequisite for enjoying the favorable $\kappa = r$ scenario associated with criteria (ii). This then leads to the following trifold distinction:

- **Balanced Regularization:** We say that a solution exhibits balanced regularization we achieve near perfect reconstruction of training samples using a minimal number of non-random/nonzero latent dimensions that convey information about \mathbf{x} .
- **Over-Regularization:** Many or all latent variables are unused or highly random such that the training data is poorly reconstructed.
- **Under-Regularization:** The reconstruction of the training data exhibits low error; however, the latent representation is redundant, i.e., an unnecessarily large number of active latent variables are used.

Achieving balanced regularization in the sense described

above does not guarantee that generated samples will be of high quality; however, as we will soon argue, over-regularized and less obviously under-regularized solutions will necessarily denigrate performance (at least when adversarial training is not used as is our focus). It therefore behooves us to pursue energy functions and/or training strategies that directly favor balanced regularization, which as will be discussed in Section 2, facilitates satisfying criteria (i). In doing so we will largely focus on a baseline VAE architecture while drawing parallels to alternative methods.

To this end, in Section 3 we describe sufficient conditions whereby the combinatorial set of distinct VAE local minima, with Gaussian encoder/decoder distributions and affine mean/variance networks, are all global optima exhibiting balanced regularization. To the best of our knowledge, this is the only such full characterization of the VAE local minima landscape involving an actual, non-trivial parameterization. Next, Section 4 details the over-regulation problem, while presenting practical strategies designed to avoid pruning away too many active dimensions at suboptimal local minima that may arise within more complex VAE (or related) architectures.

In contrast, Section 5 argues that under-regularization avoidance in a certain precise sense requires the use of energy functions with infinite gradients (or nearly so) around optima, e.g., a Gaussian VAE model with an optimal decoder variance near zero. While potentially leading to numerical issues, we describe how such energy functions are essentially unavoidable across a wide range of traditional optimization problems designed to produce maximum parsimony, as well as strategies to at least partially mitigate instabilities. We conclude with empirical examples using real-world images to corroborate our findings and demonstrate that non-adversarial training, when properly regularized, can produce reasonable samples even with full-resolution CelebA data. To the best of our knowledge, this is the first presentation of this type of realistic non-adversarial result on real-world images of this size and complexity.

2. The Value of Balanced Regularization

Over-Regularization may occur in two distinct circumstances leading to obviously deleterious effect. First, if we set $\kappa < r$, then it is generally impossible to reconstruct the training data and satisfy criteria (i).¹ Essentially all latent-variable generative modeling architectures, including GANs, will experience this straightforward form of over-regularization if κ is too small. Secondly though, if $\kappa \geq r$, then over-regularization can still potentially contribute to degraded performance. For example, the combined energy

¹This is true unless we have a decoder capacity that scales proportionally to the training data. The latter facilitates rote memorization even if $\kappa = 1 \ll r$ as discussed in (Dai et al., 2018).

function may include solutions, or possibly bad local optima, whereby more than $\kappa - r$ latent dimensions are swamped with random noise. Within a VAE this can happen if too many dimensions of $q_\phi(z|\mathbf{x})$ become independent of \mathbf{x} as observed in numerous prior works discussed in Section 4.

The impact of under-regularization as defined in Section 1 is decidedly less obvious and more contingent on the specific model architecture. But a key nuisance emerges with autoencoder structures: If $\kappa > r$ and the decoder $q_\phi(z|\mathbf{x})$ is deterministic, or at least nearly deterministic along more than r latent dimensions, then the aggregate posterior $q_\phi(z)$ will occupy zero measure within the κ -dimensional space, i.e., randomness from $\mathbf{x} \sim \mu_{gt}$ can only “fill” r dimensions. It is therefore impossible to match with a fixed prior such as $p(z) = \mathcal{N}(z|0, I)$. And even if we replace such a rigid prior with a flexible, parameterized alternative, matching degenerate distributions with near infinite density on a manifold is quite challenging, and is more-or-less just another version of our original problem, namely, matching $p_\theta(\mathbf{x})$ to μ_{gt} . It is therefore preferable to achieve balanced regularization, whereby no more or less than $\kappa - r$ dimensions of $q_\phi(z|\mathbf{x})$ produce random noise such that $q_\phi(z)$ occupies nonzero measure across κ dimensions, allowing for a more suitable match with a non-degenerate prior. (Alternatively, if $\kappa - r$ dimensions were simply set to zero, they could presumably be detected and similarly removed as well.)

In this context, balanced regularization can be viewed as a weaker prerequisite for achieving criteria (i) and (ii): it implies that we have simultaneously achieved zero reconstruction error and a latent representation such that $q_\phi(z)$ effectively occupies nonzero measure in \mathbb{R}^κ . We next examine conditions such that this is guaranteed to happen within a Gaussian VAE architecture, at least assuming we can find a local minima of the energy function

3. Sufficient Conditions for All Local Minima Exhibiting Balanced Regularization

Because our focus is continuous data, we consider a typical VAE model with a Gaussian decoder $p_\theta(\mathbf{x}|z) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$, where the mean function $\boldsymbol{\mu}_x$ is parameterized by θ and takes z as its input. As is often assumed in practice, we set $\boldsymbol{\Sigma}_x = \gamma I$, where $\gamma > 0$ is a scalar parameter within the parameter set θ . If the decoder mean function is restricted to be affine, i.e. $\boldsymbol{\mu}_x = \mathbf{W}_x z + \mathbf{b}_x$ for some weight matrix \mathbf{W}_x and bias vector \mathbf{b}_x , then the VAE energy function remains highly non-convex with potentially a combinatorial number of distinct local minima. However, we will examine conditions whereby all of these local minima are actually global optima that display balanced regularization as desired.

Although we could proceed by allowing the encoder to be arbitrarily complex, when the decoder mean function is

forced to be affine and $\Sigma_x = \gamma \mathbf{I}$, then a Gaussian encoder with an affine representation for $\mu_z = \mathbf{W}_z \mathbf{x} + \mathbf{b}_z$ and a diagonal $\Sigma_z = \text{diag}[s^2]$, where s^2 is an arbitrary non-negative parameter vector independent of \mathbf{x} , are sufficient to achieve the optimal VAE cost. Collectively, these specifications lead to the parameterization $\theta = \{\mathbf{W}_x, \mathbf{b}_x, \gamma\}$ and $\phi = \{\mathbf{W}_z, \mathbf{b}_z, s^2\}$ and energy given by $\mathcal{L}(\theta, \phi) =$

$$\int \left\{ \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\frac{1}{\gamma} \|\mathbf{x} - \mathbf{W}_x \mathbf{z} - \mathbf{b}_x\|_2^2 \right] + d \log \gamma \right. \quad (3)$$

$$\left. + \sum_{k=1}^{\kappa} (s_k^2 - \log s_k^2) + \|\mathbf{W}_z \mathbf{x} + \mathbf{b}_z\|_2^2 \right\} \mu_{gt}(d\mathbf{x}).$$

We then have the following:

Theorem 1 *For any fixed value of γ , all local minima of (3) with respect to the parameters $\{\mathbf{W}_x, \mathbf{b}_x, \mathbf{W}_z, \mathbf{b}_z, s^2\}$ are also global minima. Additionally, if the support of μ_{gt} is confined to an r -dimensional affine space, then these global minima will exhibit balanced regularization when $\gamma \rightarrow 0$, i.e., zero reconstruction error using only r active/non-random dimensions of \mathbf{z} .*

See the supplementary for proof details. This balanced regularization occurs because each possible local/global optima will necessarily define the principal subspace of the data using a minimum number of nonzero columns of \mathbf{W}_x . Furthermore, at the indices of these zero-valued columns, elements of s^2 will converge to zero while the corresponding elements of μ_z will convey information about \mathbf{x} (i.e., active, non-random dimensions). Elsewhere the variances will be set to one while the means will equal zero, indicating that no useful information about \mathbf{x} is being transferred. And because $\|\mu_z\|_0 = r$ and the reconstruction error is zero when $\gamma \rightarrow 0$, we have achieved balanced regularization per the stated definition.

We emphasize that, regardless of μ_{gt} , the number of potential minima can be combinatorially high. For example, if the covariance $\text{Cov}_{\mu_{gt}}[\mathbf{x}]$ equals $\Psi \Psi^\top$, with $\text{rank}[\Psi] = r < \kappa$, then there can exist $\frac{\kappa!}{(\kappa-r)!}$ distinct or disconnected minima of (3).² This follows from a much simpler analysis of affine VAE models from (Dai et al., 2018), where only the possibility of global (but not local) minima were considered. But since all are equally good per our results herein, this generally lessens the concern as there is no issue of suboptimal local minima with imbalanced regularization. Of course the model described by (3) is obviously a simplified version of the VAE; however, Theorem 1 nonetheless represents the only characterization of the full constellation of VAE local minima under non-trivial circumstances. It also complements the wide variety of recent efforts to analyze the complex loss surface of neural networks with linear layers (Choromanska et al., 2015a;b; Kawaguchi, 2016).

²By disconnected we mean that, to traverse from one minimum to another, we must ascend the objective function at some point along the connecting path.

4. Over-Regularization Scenarios and Practical Workarounds

As we move to more complex/deeper encoder and decoder parameterizations, it is no longer possible to characterize the complex configuration of local minima as we presented in Section 3. However, we can at least elucidate a particularly pernicious set of stationary points associated with over-regularization that may be disruptive to successful optimization trajectories in many practical settings. Note that over-regularization is well-acknowledged in the literature, with compensatory measures including heuristic annealing of the KL penalty (Bowman et al., 2015; Sønderby et al., 2016), tighter bounds on the log-likelihood (Burda et al., 2015; Rezende & Mohamed, 2015), more complex priors (Bauer & Mnih, 2018; Tomczak & Welling, 2018), or modified decoder architectures (Cai et al., 2017; Dieng et al., 2018; Yeung et al., 2017). Thus far though, published results do not indicate success generating high-resolution images, and in the majority of cases, evaluations are limited to small black-and-white and/or binary images. Additionally, unlike much existing work, in this section we analyze the problem and suggest workarounds purely from the perspective of avoiding bad VAE local minima. This increases the chances of reaching global solutions designed to achieve balanced regularization as discussed later in Section 5 and high-quality samples as shown in Section 6.

To begin, we again assume a VAE model with Gaussian encoder and decoder networks. Both encoder and decoder mean functions μ_x and μ_z , as well as the diagonal encoder covariance function $\Sigma_z = \text{diag}[\sigma_z^2]$, are computed by standard deep neural networks, with layers composed of linear weights followed by element-wise non-linear activations (the decoder covariance satisfies $\Sigma_x = \gamma \mathbf{I}$ as before). We denote the weight matrix from the first layer of the decoder mean network as $\mathbf{W}_{\mu_x}^1$, while $\mathbf{w}_{\mu_x, j}^1$ refers to the corresponding j -th column. Assuming ρ layers, we denote $\mathbf{W}_{\mu_z}^\rho$ and $\mathbf{W}_{\sigma_z^2}^\rho$ as weights from the last layers of the encoder networks producing μ_z and $\log \sigma_z^2$ respectively, with j -th rows defined as $\mathbf{w}_{\mu_z, j}^\rho$ and $\mathbf{w}_{\sigma_z^2, j}^\rho$. We then characterize the following key stationary point:

Theorem 2 *If $\mathbf{w}_{\mu_x, j}^1 = (\mathbf{w}_{\mu_z, j}^\rho)^\top = (\mathbf{w}_{\sigma_z^2, j}^\rho)^\top = \mathbf{0}$ for any $j \in \{1, 2, \dots, \kappa\}$, then the gradients of $\mathcal{L}(\theta, \phi)$ with respect to $\mathbf{w}_{\mu_x, j}^1$, $\mathbf{w}_{\mu_z, j}^\rho$, and $\mathbf{w}_{\sigma_z^2, j}^\rho$ are all equal to zero.*

Without further details, the type of stationary point described by Theorem 2 could conceivably function as a saddle point (stable or unstable), a local maximum, or a local minimum. However, in many important representative situations we would argue that it serves as the latter, or more specifically, a potentially degenerate local minimum with all or most latent dimensions suppressing information about \mathbf{x} such that an over-regularized solution is produced with high recon-

struction error. In other words, for dimensions where (??) is satisfied, the decoder will simply output i.i.d. unit-Gaussian noise that will then be zeroed out by the decoder.

For example, if we consider the VAE KL term in isolation, it follows that

$$\begin{aligned} \arg \min_{\boldsymbol{\mu}_z; \boldsymbol{\Sigma}_z \succ \mathbf{0}} \mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] &\equiv \\ \arg \min_{\boldsymbol{\mu}_z; \boldsymbol{\Sigma}_z \succ \mathbf{0}} \{ \text{tr}[\boldsymbol{\Sigma}_z] + \|\boldsymbol{\mu}_z\|_2^2 - \log|\boldsymbol{\Sigma}_z| \} &= \{ \mathbf{0}, \mathbf{I} \}. \end{aligned} \quad (4)$$

The stationary point from Theorem 2 naturally achieves this minimal solution across all j . So based on this term alone, any movement away from $\mathbf{w}_{\mu_x, j}^\rho = \mathbf{w}_{\sigma_z, j}^\rho = \mathbf{0}$ to modify $\boldsymbol{\mu}_z$ and/or $\boldsymbol{\Sigma}_z = \text{diag}[\sigma_z^2]$ will necessarily involve going against the gradient.

Of course the VAE model has a second data-dependent reconstruction factor $-\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]$. But assuming the decoder mean network is sufficiently deep, with weights near a random initialization not in the vicinity of an optimal representation, then any perturbation from $\mathbf{w}_{\mu_x, j}^1 = \mathbf{w}_{\mu_z, j}^\rho = \mathbf{w}_{\sigma_z, j}^\rho = \mathbf{0}$ may well involve an uphill battle. This is because the only way for any gradient to pass through the decoder is to perturb $\mathbf{w}_{\mu_x, j}^1$ away from zero, but the moment this happens, the large random signal coming from the encoder, which is dominated by white Gaussian noise when $\boldsymbol{\mu}_z = \mathbf{0}$ and $\boldsymbol{\Sigma}_z = \mathbf{I}$, will likely only serve to disrupt any attempted data fit after passing through the deep decoder, increasing the reconstruction cost.

Therefore, in the present aggregated circumstances, nudging $\mathbf{w}_{\mu_x, j}^1 \neq \mathbf{0}$ will be against the gradient of $-\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]$, while pushing $\boldsymbol{\mu}_z \neq \mathbf{0}$ and $\boldsymbol{\Sigma}_z \prec \mathbf{I}$ (through changes in $\mathbf{w}_{\mu_z, j}^\rho$ and $\mathbf{w}_{\sigma_z, j}^\rho$ for some/all j) to improve the predictive signal-to-noise ratio will necessarily be against the gradient of $\mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$. Hence the stationary point described by Theorem 2 can easily serve as a degenerate, over-regularized local minima. It is also important to keep in mind that arriving at such an over-regularized solution only requires the movements of the isolated weight matrices $\mathbf{W}_{\mu_x}^1$, $\mathbf{W}_{\mu_z}^\rho$ and $\mathbf{W}_{\sigma_z}^\rho$ towards zero. In contrast, to actually minimize the reconstruction error requires the complex coordination of *all* the layer weights from both the deep encoder and decoder networks, which is far more circuitous than in the affine case.

Fortunately though, multiple complementary countermeasures can be introduced to lessen the impact of this type of over-regularized local minima. To see this, consider the case where $\mathbf{W}_{\mu_x}^1 = \mathbf{W}_{\mu_z}^\rho = \mathbf{W}_{\sigma_z}^\rho = \mathbf{0}$, such that we are located in parameter space at a likely local minimum characterized by maximal over-regularization. Now suppose we add an additional affine skip connection between \mathbf{z} and \mathbf{x} , by which we mean that the decoder mean function is

modified to

$$\boldsymbol{\mu}_x = f_\theta(\mathbf{z}) + \mathbf{W}_x \mathbf{z} + \mathbf{b}_x, \quad (5)$$

where $f_\theta(\mathbf{z})$ is the original deep decoder network and \mathbf{W}_x and \mathbf{b}_x are additional parameters analogous to those defined in Section 3. If $\mathbf{W}_{\mu_x}^1 = \mathbf{0}$, then $f_\theta(\mathbf{z})$ equals some constant c that can be absorbed into \mathbf{b}_x . Hence we are now in the regime described by Theorem 1 and the analysis that follows, from which we can conclude that, even with $f_\theta(\mathbf{z})$ anchored in a bad configuration, available local minima are nonetheless obtainable with \mathbf{W}_x and \mathbf{b}_x spanning the principal subspace of the data rather than a degenerate, maximally over-regularized solution. This provides a plausible escape route. Additionally, the larger we make $\kappa = \dim[\mathbf{z}]$, the better equipped this affine bridge is to provide a preliminary fit to the data, or move further away from over-regularized local minima.

This then suggest three practical prescriptions for mitigating over-regularization:

1. Include skip connections between decoder layers,
2. Increase the value of κ such that $\kappa \gg r$ or even potentially $\kappa \approx d$ if over-regularization remains an issue.
3. If it is not practically feasible to significantly increase κ , then first train a deterministic, unregularized AE with κ set to a suitable range such that a low reconstruction error is achievable using a latent code denoted as $\tilde{\mathbf{x}}$. Then train a smaller VAE treating $\tilde{\mathbf{x}}$ as the new data and with $\tilde{d} \triangleq \dim[\tilde{\mathbf{x}}] \approx \kappa$.

With regard to prescription 2, there are multiple reasons for potentially choosing $\kappa \gg r$. First, we generally do not know r in advance, so it is obviously easier to simply choose some κ that is significantly larger rather than trying to guess $\kappa = r$. Secondly though, and equally important, the larger we make κ , *the more likely that we can fully capitalize on skip connections and avoid over-regularized solutions*. And furthermore, redundant dimensions that exist when $\kappa \gg r$ can be pruned anyway if a balanced regularization state is eventually achieved, which is possible at certain global minima as described in Section 5. Experimental results from Section 6 demonstrate that both of the above modifications are necessary in tandem.

As for optional prescription 3, an unregularized AE provides a simple way of reducing the dimensionality (and also complexity/depth) without the risk of over-regularization. A smaller/simpler VAE can then be introduced to take maximal advantage of prescriptions 1 and 2. Note that from an implementation standpoint, we adopt the reconstruction factor $\mathbb{E}_{q_\phi(\mathbf{z}|\tilde{\mathbf{x}})} \left[\frac{1}{2^\gamma} \left\| \mathbf{x} - \tilde{f}_\theta[f_\theta(\mathbf{z})] \right\|_2^2 \right]$, where \tilde{f}_θ denotes the learned AE with parameters $\tilde{\theta}$. This merely implies that

we measure the reconstruction cost in the original high-dimensional space. This is similar to methods that substitute various feature-space loss functions into the VAE energy (Hou et al., 2017).

Before proceeding, we should mention that skip connections have been proposed previously in an attempt to enhance VAE models (Cai et al., 2017; Dieng et al., 2018). However, the motivation was different, and the importance of pairing with a higher dimensional \mathbf{z} for the specific purpose of escaping bad local minima has not been explored.

5. Under-Regularization and Infinite Gradients as a Necessary Risk

We now consider necessary conditions for avoiding under-regularization in generic AE architectures, which will later be used to elucidate key VAE properties. In this regards, consider the constrained objective function $\mathcal{L}_h(\theta, \phi) =$

$$h \left(\frac{1}{dn} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - f_{\theta}(\mathbf{z}^{(i)}) \right\|_2^2 \right) + \frac{1}{d} \sum_{k=1}^{\kappa} h \left(\frac{1}{n} \|\mathbf{z}_k\|_2^2 \right)$$

s.t. $\mathbf{z}^{(i)} = g_{\phi}(\mathbf{x}^{(i)}) \quad \forall i, \theta \in \Theta,$ (6)

where $\mathbf{Z} \triangleq \{\mathbf{z}^{(i)}\}_{i=1}^n \in \mathbb{R}^{\kappa \times n}$ and \mathbf{z}_k denotes the k -th row of \mathbf{Z} . This expression can be viewed as characterizing a typical regularized AE with a generic penalty function h on the norm across training samples of each latent dimension. The multipliers $1/n$, $1/d$, and $1/(dn)$ ensure a form of proportional regularization as within energy functions composed of multiple penalty factors of varying dimension designed to favor sparsity (Wipf & Wu, 2012). The square-root Lasso can be viewed as a special case of this strategy that emerges when h is a square-root function (Belloni et al., 2011). We adopt this formalism to avoid distracting complications from tunable trade-off parameters; however, our central conclusions still hold even when such a parameter is introduced. And finally, the constraint $\theta \in \Theta$ is included to prevent the trivial solution $\mathbf{Z} \rightarrow \mathbf{0}$, which could occur if each $\mathbf{z}^{(i)}$ is pushed to zero while f_{θ} includes an unconstrained compensatory factor that grows towards infinity such that the error $\|\mathbf{x}^{(i)} - f_{\theta}(\mathbf{z}^{(i)})\|_2$ can still be minimized to zero. Any regularized AE must include such constraints to avoid trivial solutions, or else additional penalty terms on θ and ϕ that serve a similar purpose.

Given a generic AE architecture as in (6), it is natural to examine what possible functions h are such that any global minimum of $\mathcal{L}_h(\theta, \phi)$ is guaranteed to exhibit balanced regularization. This can be addressed as follows:

Theorem 3 Assume the constraint $\theta \in \Theta$ and data $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^n \in \mathbb{R}^{d \times n}$ are such that to achieve $\mathbf{x}^{(i)} = f_{\theta}(\mathbf{z}^{(i)})$ for all i (i.e., perfect reconstruction) requires that $\|\mathbf{z}_k\|_2 > 0$ for at least $r < d$ rows of \mathbf{Z} . Then to guarantee that

minimization of $\mathcal{L}_h(\theta, \phi)$ achieves zero reconstruction error using at most r nonzero rows of \mathbf{Z} (i.e., active dimensions), h must have an unbounded gradient around zero.

Note that a similar result can be obtained by replacing the reconstruction penalty with the additional constraint $\sum_{i=1}^n \|\mathbf{x}^{(i)} - f_{\theta}(\mathbf{z}^{(i)})\|_2^2 = 0$, in which case no trade-off parameter, fixed or otherwise, need be included. We also emphasize that Theorem 3 effectively implies that, to guarantee every global optima displays balanced regularization per our definition, the constituent penalty functions must have an infinite gradient around zero. Given that we may readily introduce arbitrary scalings and translations, this condition is more-or-less tantamount to requiring penalty functions with an energy gap that is unbounded about zero. For example, the selection $h(u) = \mathcal{I}[u > 0]$, i.e., an indicator function that equals zero if $u = 0$ and one for all $u > 0$, will guarantee that any global minimum of (6) displays balanced regularization under the stated conditions. However, given that $h(u) \equiv \lim_{p \rightarrow 0} u^p$ and $\lim_{p \rightarrow 0} \frac{1}{p}(u^p - 1) = \log u$, we see that an unbounded log function can essentially achieve the same result in the limit.

The VAE can be interpreted as a form of stochastic AE, with subtle regularization effects introduced via the interplay between the reconstruction and KL terms. A number of recent works have mentioned that if a flexible decoder variance parameter γ is included within a Gaussian VAE, then the optimal value may converge to zero, resulting in infinite gradients and potential instabilities (Dai & Wipf, 2018; Mattei & Frelsen, 2018; Takahashi et al., 2018). While unbounded gradients may indeed be troublesome from an optimization perspective, based on the analysis of this section, we frame such gradients as a necessary component of any model that attempts to produce balanced regularization. In this regard, it has been argued that as $\gamma \rightarrow 0$, the VAE can achieve zero reconstruction error at the global optimum by selectively pushing $q_{\phi}(\mathbf{z}|\mathbf{x})$ towards a degenerate Gaussian, with zero variance along the minimal number of directions needed for reconstructing \mathbf{x} , and unit variance elsewhere so as to minimize the increase in the KL regularization factor (Dai & Wipf, 2018). This is exactly a stochastic version of balanced regularization, and ensures that $q_{\phi}(\mathbf{z})$ occupies full measure in κ -dimensional space, even if $q_{\phi}(\mathbf{z}) \neq p(\mathbf{z})$.

Furthermore, because $q_{\phi}(\mathbf{z})$ is not a troublesome degenerative distribution, it is ripe for pairing with any standard generative modeling paradigm that facilitates tractable sampling. For example, we could apply a second-stage VAE as proposed in (Dai & Wipf, 2018), or invertible networks capable of warping a κ -dimensional Gaussian into an arbitrary κ -dimensional distribution with full measure (Dinh et al., 2016). Regardless, we contend that achieving balanced regularization represents an essential stepping stone, even if infinite gradients constitute a necessary evil.

We close this section by acknowledging that energy functions involving infinite gradients and/or unbounded regions are already indispensable across a wide range of sparse estimation problems and structured regression (Gorodnitsky & Rao, 1997). This history implies that when training a VAE or other related autoencoder structure, we may borrow appropriate tools designed to mitigate the risk of converging to bad local solutions or regions of instability. In this vein, one effective strategy involves partially minimizing what amounts to a smoothed version of the original objective function. The degree of smoothness is then gradually reduced as the optimization trajectory moves towards an optimum. This procedure, which serves as a form of homotopy continuation method, is frequently used to find maximally sparse representations with minimal reconstruction error (Chartrand & Yin, 2008; Hu et al., 2012; Xu et al., 2013).

The VAE accomplishes something similar when we choose to iteratively estimate γ during training rather than merely setting its value to near zero as may be theoretically optimal. Initially, when the reconstruction cost is still high, γ will be relatively large and the overall VAE energy will be relatively smooth. It is only later as $\sum_{i=1}^n \|\mathbf{x}^{(i)} - f_{\theta}(\mathbf{z}^{(i)})\|_2^2$ becomes small that γ will follow suite, and by this point it is more likely that we have already approached a basin of attraction with balanced regularization.

6. Experimental Validation and Discussion

A Simple Over- and Under-Regularization Test: We first perform experiments to examine robustness to over-regularization using the prescriptions from Section 4 and under-regularization using a Gaussian decoder with trainable variance parameter γ per Section 5. The objective at this point is to achieve optimal reconstructions using minimal latent dimensions (i.e., balanced regularization), as opposed to the final goal of generating high-quality samples which will be addressed later. We train a VAE model on MNIST data (LeCun et al., 1998). The encoder is a 6-layer feed-forward network and the decoder is a residual network (He et al., 2016) with 3 residual blocks, each containing 2 layers, to instantiate skip connections. Note that a Gaussian decoder as we have adopted is unlike many previous works that simplify the problem by binarizing the data and then applying a Bernoulli decoder appropriate for binary images. We vary the latent dimension κ from 50 to 784, where $d = 784$ is the dimension of the data itself. See the supplementary file for training details. Three different metrics are adopted to evaluate models: the reconstruction error, the negative log-likelihood (NLL) on the test set, and the number of active latent dimensions per the metric from (Burda et al., 2015).

Table 1 displays the results, where we observe that as κ becomes sufficiently large, the reconstruction error mostly decreases while the number of active dimensions saturates

κ	Rec. Err.	NLL	# Active
50	4.98	-315	10
100	3.42	-580	15
200	2.74	-718	20
400	2.57	-743	22
784	2.78	-806	22

Table 1. Reconstruction error (test), NLL (test), and number of active dimensions on MNIST as κ is varied.

	CelebA-64		CelebA-128	
	Rec. Err.	MMD	Rec. Err.	MMD
Learnable γ	60.5	20.24	352.8	93.3
Fix $\gamma = \gamma^*$	59.6	69.59	349.9	291.8

Table 2. Reconstruction error and MMD between $q_{\phi}(\mathbf{z})$ and $\mathcal{N}(0, \mathbf{I})$. We first train a VAE with learnable γ and obtain the optimal value γ^* . Then we fix $\gamma = \gamma^*$ and re-train the same network from scratch. Though the final reconstruction errors are almost the same, the MMDs between $q_{\phi}(\mathbf{z})$ and the standard $\mathcal{N}(0, \mathbf{I})$ are significantly different.

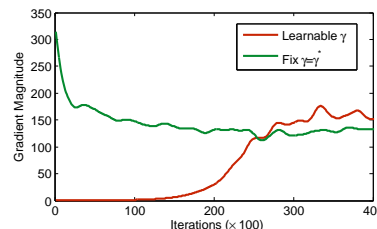


Figure 1. The Evolution of the gradient $\left\| \frac{d\mathcal{L}(\theta, \phi)}{d\mathbf{z}} \right\|_2$. Although both curves end up with similar final values, the large initial gradient with fixed γ is disruptive to the final solution.

around 20. Overall the NLL on the test set improves even up to $\kappa = 784$, indicating that the model has likely not under-regularized (or overfit) as well, presumably because excessive active dimensions have been pruned. Interestingly though, if we remove the skip connections, this relatively deep model completely over-regularizes, and all the active dimensions are turned off leading to completely useless reconstructions (not shown). So indeed both the skip connections and sufficiently large κ are needed.

Under-Regularization and the Benefits of Learning γ :

Although a small γ may lead to large VAE gradients, we have argued that this may constitute a necessary risk if we want to obtain minimal reconstruction error using the fewest number of active latent dimensions, i.e. avoid under-regularization. We now demonstrate that learning γ , as a form of homotopy continuation method, is better than fixing it to an optimal small value as discussed in Section 5. In particular, we first train a VAE on the CelebA dataset (Liu et al., 2015), both 64×64 and 128×128 resolution versions, and learn an appropriate small value of γ denoted γ^* (see supplementary file for network and training details). We then retrain the same network from scratch but with $\gamma = \gamma^*$ fixed. The resulting models are evaluated via the recon-

struction error and the maximum mean discrepancy (MMD) between the aggregated posterior $q_\phi(\mathbf{z})$ and the prior $p(\mathbf{z})$. If too few latent dimensions are removed by swamping the appropriate channels with noise following the prior, then we would expect $q_\phi(\mathbf{z})$ to be confined to a manifold in \mathbb{R}^κ and the MMD to be much larger.

Results are displayed in Table 2, where as expected, the reconstruction errors are nearly identical, but the learnable γ has much lower MMD values. For the 128×128 case, we also plot the evolution of the gradient magnitudes $\left\| \frac{d\mathcal{L}(\theta, \phi)}{d\mathbf{z}} \right\|_2$ in Figure 1. When γ is learned (red curve), the gradient increases slowly; however, with fixed $\gamma = \gamma^*$, there exists a huge gradient right from the start since γ^* is small but the reconstruction error is high. This contributes to a worse final solution per the Table 2 results.

From Balanced Regularization to Good Samples: Thus far we have merely investigated strategies for achieving balanced regularization. We now apply these ideas to the original goal of generating samples of high-resolution images. Note that the vast majority of AE-structured, non-adversarial approaches are only tested using NLL scores or related; however, this is largely orthogonal to the goal of improving quantitative measures of visual quality (Theis et al., 2016). Instead, metrics such as the Fréchet Inception Distance (FID) (Heusel et al., 2017) have been designed for this purpose, as well as complementary scores like the recently proposed Number of statistically-Different Bins (NDB) (Richardson & Weiss, 2018). The latter evaluates generation performance in the original images space by binning samples into K learned clusters; it therefore does not depend on subtleties of inception network like FID, and provides information about high-level sample diversity.

We quantitatively evaluate networks on FID and NDB scores using the original, full-resolution CelebA data center-cropped to 128×128 . Based on the analysis herein, we first train an unregularized AE to project the original data down to a $\tilde{d} = 64$ -dimensional feature space. In this regime, it is feasible to train a small VAE with $\tilde{d} = \kappa$ and skip connections to achieve balanced regularization (at least assuming that \tilde{d} is sufficiently large). And finally, per the suggestion of (Dai & Wipf, 2018), we train a second small VAE to sample from the aggregate posterior $q_\phi(\mathbf{z})$ produced by the first, which should benefit from balanced regularization. Network and training details are in the supplementary. We refer to this model as VAE+, which is simple to train since nearly all of the parameters reside in the AE structure.

Unlike other non-adversarial AE-based models, the WAE-MMD approach (Tolstikhin et al., 2018) has also been quantitatively evaluated on real-world color images like CelebA, but only on a downsampled version. For balanced comparison, we train a WAE-MMD network using the aggregate AE plus VAE structure described above. As an additional base-

	Recon. FID	Gen. FID
Downsample	39.2	39.2
WAE-MMD	59.8	66.9
VAE	45.9	64.8
VAE+	45.6	47.5

Table 3. Reconstruction and generation FID scores on 128×128 CelebA dataset (smaller is better). *Downsample* refers to comparing the original images with downsampled 64×64 versions. This seemingly innocuous transformation, which preserves all the major structures, can produce an FID almost as high as VAE+. This indicates that the FID may be sensitive to minor details, which could impact the performance of AE-structured models.

	$K = 100$	$K = 200$	$K = 300$
Downsample	0.05	0.06	0.04
WAE-MMD	0.64	0.62	0.58
VAE	0.88	0.85	0.81
VAE+	0.29	0.28	0.24

Table 4. NDB scores with varying numbers of bins (K). *Downsample* is defined as in Table 3; this value is near 0 as expected.

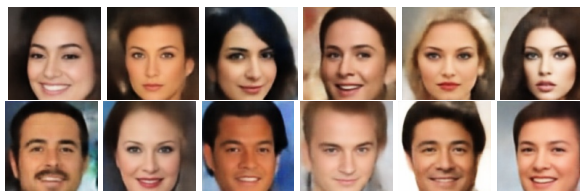


Figure 2. Select generated samples

line, we repeat this process with a plain VAE model using the same architecture. The reconstruction and generation FID scores are shown in Table 3. The reconstruction FID, which is only sensitive to criteria (i), can be regarded as a lower bound of the generation FID score since it is unlikely that the generated samples could ever be better than the reconstructed ones. The fact that the VAE+ generation FID is so close to the reconstruction FID indicates that it has done well in addressing criteria (ii) via its sensitivity to balanced regularization. We also report results in Table 4 using the complementary NDB metric. The VAE+ again performs best by a significant margin, indicating a higher level of sample diversity. And finally, we show some of the good samples generated by VAE+ in Figure 2. More samples from VAE+ and other methods are in the supplementary.

Discussion: We have analyzed the regularization balance of AE-structured models in general and VAEs in particular, including local minima properties and necessary energy function characteristics. This leads to useful practical prescriptions and the first demonstration of high-quality, diverse generation results from AE-structured, non-adversarial training on 128×128 color images. The only other attempt at this resolution comes from (Cai et al., 2017), but the generated images are comparably blurry (see supplementary) and no quantitative evaluation has been demonstrated.

References

- 440 Arora, S. and Zhang, Y. Do gans actually learn the
441 distribution? an empirical study. *arXiv preprint*
442 *arXiv:1706.08224*, 2017.
443
444
- 445 Bauer, M. and Mnih, A. Resampled priors for variational
446 autoencoders. *arXiv preprint arXiv:1810.11428*, 2018.
447
- 448 Belloni, A., Chernozhukov, V., and Wang, L. Square-root
449 lasso: pivotal recovery of sparse signals via conic pro-
450 gramming. *Biometrika*, 98(4):791–806, 2011.
- 451 Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A. M., Joze-
452 fowicz, R., and Bengio, S. Generating sentences from
453 a continuous space. *arXiv preprint arXiv:1511.06349*,
454 2015.
455
- 456 Burda, Y., Grosse, R., and Salakhutdinov, R. Importance
457 weighted autoencoders. *arXiv preprint arXiv:1509.00519*,
458 2015.
- 459 Cai, L., Gao, H., and Ji, S. Multi-stage variational auto-
460 encoders for coarse-to-fine image generation. *arXiv*
461 *preprint arXiv:1705.07202*, 2017.
- 462
- 463 Chartrand, R. and Yin, W. Iteratively reweighted algorithms
464 for compressive sensing. *Proc. Int. Conf. Acoustics,*
465 *Speech, and Signal Proc.*, 2008.
466
- 467 Choromanska, A., Henaff, M., Mathieu, M., Arous, G., and
468 LeCun, Y. The loss surfaces of multilayer networks. In
469 *International Conference on Artificial Intelligence and*
470 *Statistics*, 2015a.
- 471 Choromanska, A., LeCun, Y., and Arous, G. Open problem:
472 The landscape of the loss surfaces of multilayer networks.
473 In *Conference on Learning Theory*, 2015b.
474
- 475 Dai, B. and Wipf, D. Diagnosing and enhancing gaussian
476 VAE models. *Advances in Neural Information Processing*
477 *Systems, Bayesian Deep Learning Workshop*, 2018.
- 478
- 479 Dai, B., Wang, Y., Aston, J., Hua, G., and Wipf, D. Connec-
480 tions with robust PCA and the role of emergent sparsity
481 in variational autoencoder models. *Journal of Machine*
482 *Learning Research*, 2018.
- 483 Dieng, A. B., Kim, Y., Rush, A. M., and Blei, D. M. Avoid-
484 ing latent variable collapse with generative skip models.
485 *arXiv preprint arXiv:1807.04863*, 2018.
486
- 487 Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density esti-
488 mation using real nvp. *arXiv preprint arXiv:1605.08803*,
489 2016.
- 490 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B.,
491 Warde-Farley, D., Ozair, S., Courville, A., and Bengio,
492 Y. Generative adversarial networks. In *arXiv preprint*
493 *arXiv:1406.2661*, 2014.
494
- Gorodnitsky, I. F. and Rao, B. D. Sparse signal recon-
struction from limited data using focuss: A re-weighted
minimum norm algorithm. *IEEE Transactions on signal*
processing, 45(3):600–616, 1997.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings
in deep residual networks. In *European conference on*
computer vision, pp. 630–645. Springer, 2016.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and
Hochreiter, S. GANs trained by a two time-scale update
rule converge to a local Nash equilibrium. In *Advances in*
Neural Information Processing Systems, pp. 6626–6637,
2017.
- Hou, X., Shen, L., Sun, K., and Qiu, G. Deep feature
consistent variational autoencoder. In *Applications of*
Computer Vision, 2017 IEEE Winter Conference on, pp.
1133–1141. IEEE, 2017.
- Hu, Y., Lingala, S. G., and Jacob, M. A fast majorize-
minimize algorithm for the recovery of sparse and low-
rank matrices. *IEEE Transactions on Image Processing*,
21(2):742–753, 2012.
- Kawaguchi, K. Deep learning without poor local minima.
In *Advances in Neural Information Processing Systems*,
2016.
- Kingma, D. and Welling, M. Auto-encoding variational
Bayes. In *International Conference on Learning Repre-*
sentations, 2014.
- Kingma, D., Salimans, T., Jozefowicz, R., Chen, X.,
Sutskever, I., and Welling, M. Improved variational in-
ference with inverse autoregressive flow. In *Advances in*
Neural Information Processing Systems, pp. 4743–4751,
2016.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-
based learning applied to document recognition. *Proceed-*
ings of the IEEE, 86(11):2278–2324, 1998.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face
attributes in the wild. In *IEEE International Conference*
on Computer Vision, pp. 3730–3738, 2015.
- Lucic, M., Kurach, K., Michalski, M., Gelly, S., and Bous-
quet, O. Are GANs created equal? A large-scale study.
Advances in Neural Information Processing Systems,
2018.
- Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., and
Frey, B. Adversarial autoencoders. *arXiv preprint*
arXiv:1511.05644, 2016.
- Mattei, P.-A. and Frellsen, J. Leveraging the exact like-
lihood of deep latent variables models. *arXiv preprint*
arXiv:1802.04826, 2018.

- 495 Metz, L., Poole, B., Pfau, D., and Sohl-Dickstein, J. Un-
496 rolled generative adversarial networks. *arXiv preprint*
497 *arXiv:1611.02163*, 2016.
- 498 Rezende, D., Mohamed, S., and Wierstra, D. Stochastic
499 backpropagation and approximate inference in deep gen-
500 erative models. In *International Conference on Machine*
501 *Learning*, 2014.
- 503 Rezende, D. J. and Mohamed, S. Variational inference with
504 normalizing flows. *arXiv preprint arXiv:1505.05770*,
505 2015.
- 507 Richardson, E. and Weiss, Y. On gans and gmms. *arXiv*
508 *preprint arXiv:1805.12462*, 2018.
- 509 Sønderby, C. K., Raiko, T., Maaløe, L., Sønderby, S. K.,
510 and Winther, O. How to train deep variational autoen-
511 coders and probabilistic ladder networks. *arXiv preprint*
512 *arXiv:1602.02282*, 2016.
- 514 Takahashi, H., Iwata, T., Yamanaka, Y., Yamada, M., and
515 Yagi, S. Student-t variational autoencoder for robust
516 density estimation. In *International Joint Conference on*
517 *Artificial Intelligence*, pp. 2696–2702, 2018.
- 519 Theis, L., van den Oord, A., and Bethge, M. A note on
520 the evaluation of generative models. In *International*
521 *Conference on Learning Representations*, pp. 1–10, 2016.
- 522 Tolstikhin, I., Bousquet, O., Gelly, S., and Schoelkopf, B.
523 Wasserstein auto-encoders. *International Conference on*
524 *Learning Representations*, 2018.
- 526 Tomczak, J. and Welling, M. VAE with a VampPrior. In
527 *International Conference on Artificial Intelligence and*
528 *Statistics*, pp. 1214–1223, 2018.
- 530 van den Berg, R., Hasenclever, L., Tomczak, J. M., and
531 Welling, M. Sylvester normalizing flows for variational
532 inference. In *Uncertainty in Artificial Intelligence*, 2018.
- 533 Wipf, D. and Wu, Y. Dual-space analysis of the sparse linear
534 model. In *Advances in Neural Information Processing*
535 *Systems*, pp. 1745–1753, 2012.
- 537 Xu, L., Zheng, S., and Jia, J. Unnatural l0 sparse represen-
538 tation for natural image deblurring. In *IEEE Conference*
539 *on Computer Vision and Pattern Recognition*, pp. 1107–
540 1114, 2013.
- 542 Yeung, S., Kannan, A., Dauphin, Y., and Fei-Fei, L. Tack-
543 ling over-pruning in variational autoencoders. *arXiv*
544 *preprint arXiv:1706.03643*, 2017.
- 545
546
547
548
549