# A Spectral Regularizer for Unsupervised Disentanglement

**Aditya Ramesh** [1 2]   **Youngduck Choi** [3]   **Yann LeCun** [4]

## Abstract

A generative model with a disentangled representation allows for control over independent aspects of the output. Learning disentangled representations has been a recent topic of great interest, but it remains poorly understood. We show that even for GANs that do not possess disentangled representations, one can find curved trajectories in latent space over which local disentanglement occurs. These trajectories are found by iteratively following the leading right-singular vectors of the Jacobian of the generator with respect to its input. Based on this insight, we describe an efficient regularizer that aligns these vectors with the coordinate axes, and show that it can be used to induce disentangled representations in GANs, in a completely unsupervised manner.

## 1. Introduction

Rapid progress has been made in the development of generative models capable of producing realistic samples from distributions over complex, high-resolution natural images (Karras et al. (2017), Brock et al. (2018)). Despite the success of these models, it is unclear how one can achieve control over the data-generation process when labeled instances are not available. Perturbing an individual component of a latent variable typically results in an unpredictable change to the output. When the model has a disentangled representation, these changes become interpretable, and each component of the latent variable affects a distinct attribute. So far, the problem of learning disentangled reprersentations remains poorly understood. Most approaches for doing this have focused on VAEs (Kingma & Welling (2013), Rezende et al. (2014)), which produce blurry samples in practice. The few approaches that have been developed for GANs, such as InfoGAN (Chen et al., 2016), have had comparably limited

[1]OpenAI [2]Continuation of work done at New York University [3]Department of Computer Science, Yale University [4]Department of Computer Science, New York University. Correspondence to: Aditya Ramesh <aramesh@openai.com>.
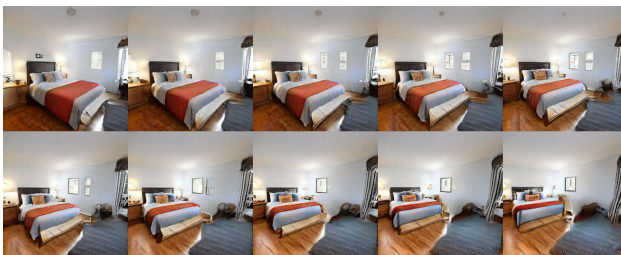
Figure 1: Trajectory obtained by following the leading eigenvector of $M_z(z_0)$ from a fixed embedding $z_0 \in \mathbb{R}^{256}$ in the latent space of a GAN trained on the LSUN Bedrooms dataset at $256 \times 256$ resolution. The camera angle is varied smoothly, while the content of the bedroom is preserved. The parameters used for Algorithm 1 are given by $\alpha, \rho, N = 1.5 \cdot 10^{-2}, 0.99, 2000$; we show iterates $0, 20, \ldots, 180$. The architecture and training procedure are as described in (Mescheder et al., 2018). Videos of the trajectories for the first four eigenvectors can be viewed at this URL: https://drive.google.com/open?id=1_SxlvcjakzkL8vNY4hporwtvg8mLQw4S.

success.

Learning disentangled representations is of interest, because they allow us gauge the ability of generative models to benefit downstream tasks. For example, a robotics application may require a vision model that reliably determines the orientation of an object, subject to changes in viewpoint, lighting conditions, and intra-class variation in appearance. A generative model with independent control over the full extent of each of these attributes may help to achieve the level of robustness that is required for these applications. Moreover, the ability to synthesize concepts in ways that are not present in the original data (e.g., rotation of an object to a position not encoutered during training) is a useful benchmark for reasoning and out-of-distribution generalization. Often times, the goal of training a generative model is to observe that this type of generalization occurs.

Following (Higgins et al., 2016), we say a generative model $G : Z \to X$, where $Z \subset \mathbb{R}^m$ and $X \subset \mathbb{R}^n$, possesses a disentangled representation when it satisfies two important properties. The first property is *independence* of the components of the latent variable. Roughly speaking, this means that a perturbation to any component of a latent variable should result in a change to the output that is distinct, in some sense, from the changes resulting from perturbations to the other components. For example, if the first component of the latent variable controls hair length, then the other components should not influence this attribute. The

second property is *interpretability* of the changes in $X$ resulting from these perturbations. Suppose that the process of sampling from the distribution $p_{\text{data}}$ modeled by $G$ can be realized by a simulator that is parameterized over a list of independent, scalar-valued factors or attributes. These attributes might correspond to concepts such as lighting, azimuth, gender, age, and so on. This property is met when $G(z + \alpha e_i)$ results in a change along exactly one of these attributes, for each $i \in [1, m]$, where $e_i$ is the $i$th standard basis vector and $\alpha > 0$. Quantifying the extent to which this property holds is generally challenging.

Classical dimensionality reduction techniques, such as PCA and metric MDS, can be used to obtain a latent representation satisfying the first property. For instance, the former ensures that changes to the output resulting from changes to distinct components of the latent variable are orthogonal, which is a particularly restrictive form of distinctness. These techniques are guaranteed to faithfully represent the underlying structure only when the data occupy a linear subspace. For many applications of interest, such as image modeling, the data manifold could be a highly curved and twisted surface that does not satisfy this assumption. The Swiss Roll dataset is a simple example for which this is the case: two points with small euclidean distance could have large geodesic distance on the data manifold. Since both techniques measure similarity between each pair of points using euclidean distance, distant outputs in $X$ could get mapped to nearby embeddings in $Z$. Thus, the second property often fails to hold in practice.

Deep generative models have the potential to learn representations that satisfy both of the required properties for disentanglement. Our work focuses on the case in which $G$ is a GAN generator. Let $J_G(z_0) \in \mathbb{R}^{n \times m}$ be the Jacobian of $G$ evaluated at the latent variable $z_0 \in Z$. Our main contributions are as follows:

1. We show that the leading right-singular vectors of $J_G(z_0)$ can be used to obtain a local disentangled representation about a neighborhood of $z_0$;

2. We show that following the path determined by a leading right-singular vector in the latent space of a GAN yields a trajectory along which local changes to the output are interpretable; and

3. We formulate a regularizer that induces disentangled representations in GANs, by aligning the top $k$ right-singular vectors $v_1, \ldots, v_k$ with the first $k$ coordinate directions, $e_1, \ldots, e_k$.

## 2. Related Work

To date, the two most successful approaches for unsupervised learning of disentangled representations are the $\beta$-VAE (Higgins et al., 2016) and InfoGAN (Chen et al., 2016). The former proposes increasing the weight $\beta$ for the KL-divergence term in the objective of the VAE (Kingma & Welling (2013), Rezende et al. (2014)), which is normally set to one. This weight controls the tradeoff between minimizing reconstruction error, and minimizing the KL-divergence between the approximate posterior and the prior of the decoder. The authors observe that as the KL-divergence is reduced, more coordinate directions in the latent space of the decoder end up corresponding to disentangled factors. To measure the extent to which this occurs, they describe a disentanglement metric score, which we use in Section 6. Follow-up work (Burgess et al., 2018) analyzes $\beta$-VAE from the perspective of information channel capacity, and describes principled way of controlling the tradeoff between reconstruction error and disentanglement. Several variants of the $\beta$-VAE have also since been developed (Kim & Mnih (2018), Chen et al. (2018), Esmaeili et al. (2018)).

InfoGAN Chen et al. (2016) augments the GAN objective (Goodfellow et al., 2014) with a term that maximizes the mutual information between the generated samples and small subset of latent variables. This is done by means of an auxiliary classifier that is trained along with the generator and the discriminator. Adversarial training offers the potential for a high degree of realism that is difficult to achieve with models like VAEs, which are based on reconstruction error. However, InfoGAN finds fewer disentangled factors than VAE-based approaches on datasets such as CelebA (Liu et al., 2015) and 3DChairs (Aubry et al., 2014), and offers limited control over each latent factor (e.g., maximum rotation angle for azimuth). The mutual information regularizer also detriments sample quality. Hence, the development of an unsupervised method for learning disentangled representations with GANs, while meeting or exceeding the quality of those found by VAE-based approaches, remains an open problem.

Our work makes an important step in this direction. In contrast to previous approaches, which are based on information theory, we leverage the spectral properties of the Jacobian of the generator. This new perspective not only allows us to induce high-quality disentangled representations in GANs, but also to find curved paths over which disentanglement occurs, for GANs that do not possess disentangled representations.

## 3. Identifying Local Disentangled Factors

We begin by describing how the singular value decomposition of $J_G(z_0)$ can be used to define a local generative model about $z_0$ that satisfies the first property of disentangled representations. In particular, we show that the left-singular vectors of $J_G(z_0)$ form an orthonormal set of directions from $G(z_0)$ along which the magnitudes of the instantaneous changes in $G$ are maximized. Since perturbations

from $z_0$ along the right-singular vectors result in changes from $G(z_0)$ along the left-singular vectors, the right-singular vectors result in distinct changes to the output of $G$. This simple relationship allows us to define a local generative model that satisfies the independence property.

First, we show that the perturbations from $z_0$ along the right-singular vectors maximize the magnitude of the instantaneous change in $G(z_0)$. Given an arbitrary vector $v \in \mathbb{R}^m$, the directional derivative

$$\lim_{\epsilon \to 0} \frac{G(z_0 + \epsilon v) - G(z_0)}{\epsilon} = J_G(z_0)v. \qquad (1)$$

measures the instantaneous change in $G$ resulting from a perturbation along $v$ from $z_0$. The magnitude of this change is given by

$$\|J_G(z_0)v\| = (v^t J_G(z_0)^t J_G(z_0)v)^{1/2}$$
$$=: (v^t M_z(z_0)v)^{1/2} =: n_z(v, z_0), \qquad (2)$$

which is a seminorm involving the positive semidefinite matrix $M_z(z_0) := J_G(z_0)^t J_G(z_0) \in \mathbb{R}^{m \times m}$. The unit-norm perturbation from $z_0$ that maximizes $n_z(v, z_0)$ is given by

$$v_1 := \max_{v \in \mathbb{S}^{m-1}} n_z(v, z_0) = \max_{v \in \mathbb{S}^{m-1}} v^t J_G(z_0)^t J_G(z_0)v, \qquad (3)$$

where $\mathbb{S}^{m-1} \subset \mathbb{R}^m$ is the unit sphere. This is the first eigenvector of $M_z(z_0)$, which coincides with the first right-singular vector of $J_G(z_0)$. It follows from the singular value decomposition of $J_G(z_0)$ that the first left-singular vector is given by $u_1 = \sigma_1^{-1} J_G(z_0)v_1$, where $\sigma_1$ is the first singular value. Hence, a perturbation from $z_0$ along $v_1$ maximizes the magnitude of the instantaneous change in $G(z_0)$, and this change occurs along $u_1$.

Next, we consider the unit-norm perturbation orthogonal to $v_1$ that maximizes $n_z(v, z_0)$. It is given by

$$v_2 := \max_{v \in \mathbb{S}^{d-1} \cap \text{span}(v_1)^\perp} v^t J_G(z_0)^t J_G(z_0)v. \qquad (4)$$

This is the second eigenvector of $M_z(v, z_0)$, which coincides with the second right-singular vector of $J_G(z_0)$. As before, we get $u_2 = \sigma_2^{-1} J_G(z_0)v_2$, where $\sigma_2$ is the second singular value. So a perturbation from $z_0$ along $v_2$ results in an instantaneous change in $G(z_0)$ along $u_2$. Continuing in this way, we consider the $k$th unit-norm perturbation orthogonal to $v_1, \ldots, v_{k-1}$ that maximizes $n_z(v, z_0)$, for each $k \in [2, r]$, where $r := \text{rank}(M_z(z_0))$. This shows that the right-singular vectors of $J_G(z_0)$ maximize the magnitude of the instantaneous change in magnitude of $G(z_0)$, and these changes occur along the corresponding left-singular vectors.

Now, we use the right-singular vectors of $J_G(z_0)$ to define a local generative model about $z_0$. Consider the func-
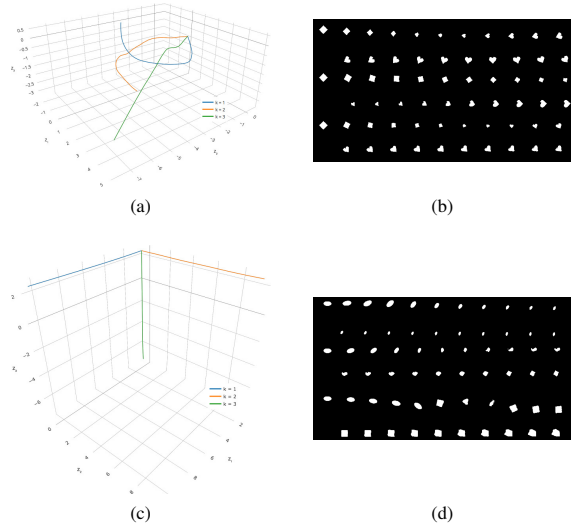


(a)     (b)



(c)     (d)

Figure 2: Trajectories obtained by following the first three eigenvectors of $M_z(z_0)$ from a fixed embedding $z_0 \in \mathbb{R}^3$ in the latent spaces of two GANs with identical architectures, trained on the dSprites dataset (Matthey et al., 2017). Subfigure (a) shows plots of the trajectories $\{\gamma_k\}_{k \in [1,3]}$, and subfigure (b) shows the outputs of $G$ at iterates $0, 100, \ldots, 2000$ of each trajectory $\gamma_k$, for $k \in [1, 3]$, from top to bottom, respectively. Subfigures (c) and (d) show the same information for an alignment-regularized GAN ($k, \lambda = 3, 0.1$); the trajectories are now axis-aligned, as expected. The parameters used for Algorithm 1 are given by $\alpha, \rho, N = 5 \cdot 10^{-3}, 0.99, 2000$. Details regarding the model architecture are given in Appendix C.

tion $\bar{G}_{z_0} : \mathbb{R}^r \to \mathbb{R}^n$,

$$\bar{G}_{z_0} : \alpha \mapsto G\left(z_0 + \sum_{i \in [1,r]} \alpha_i v_i\right).$$

The components of $\alpha$ control perturbations along orthonormal directions, and these directions also result in orthonormal changes to $G(z_0)$. Hence, $\bar{G}_{z_0}$ satisfies the first property for a generative model to possess a disentangled representation, but only about a neighborhood of $z_0$. Figure 13 in Appendix F investigates whether $\bar{G}_{z_0}$ also satisfies the second property: interpretability of changes to individual components of $\alpha$. We can see that perturbations along the leading eigenvectors $M_z(z_0)$, especially the principal eigenvector, often result in the most drastic changes. These changes are interpretable, and tend to make modifications to isolated attributes of the face. To see this in more detail, we consider the top two rows of subfigure (g). Movement along the first two eigenvectors changes hair length and facial orientation; movement along the third eigenvector decreases the length of the bangs; movement along the fourth and fifth eigenvectors changes background color; and movement along the sixth and seventh eigenvectors changes hair color.

## 4. Finding Quasi-Disentangled Paths

Generative models known to possess disentangled representations, such as $\beta$-VAEs (Higgins et al., 2016), allow for

continuous manipulation of attributes via perturbations to individual coordinates. Starting from a latent variable $z_0 \in Z$, we can move along the path $\gamma : t \mapsto z_0 + te_i$ in order to vary a single attribute of $G(z_0)$, while keeping the others held fixed. GANs are not known to learn disentangled representations automatically. Nonetheless, the previous section shows that the local generative model $\bar{G}_{z_0}$ does possess a disentangled representation, but only about a neighborhood of the base point $z_0$. We explore whether it is possible to extend this local model to obtain disentanglement along a continuous path from $z_0$. To do this, we construct a trajectory $\gamma_k : \mathbb{R} \to \mathbb{R}^n$, $t \mapsto G(\gamma_k(t))$ by repeatedly following the $k$th leading eigenvector of $M_z(\gamma_k(t))$, where $\gamma(0) := z_0$. The procedure used to do this is given by Algorithm 1.

We first test the procedure on a toy example for which it is possible to explicitly plot the trajectories. We use the dSprites dataset (Matthey et al., 2017), which consists of $64 \times 64$ white shapes on black backgrounds. Each shape is completely determined by six attributes: symbol (square, ellipse, or heart), scale (6 values), rotation angle from $0$ to $2\pi$ (40 values), and $x$- and $y$-positions (30 values each). We trained a GAN on this dataset using a latent variable size of three, fewer than the six latent factors that determine each shape. Figure 2 shows that outputs of the generator along these trajectories vary locally along only one or two attributes at a time. Along the first trajectory $\gamma_1$, the generator first decreases the scale of the square, then morphs it into a heart, increases the scale again, and finally begins to rotate it. Similar comments apply to the other two trajectories, $\gamma_2$ and $\gamma_3$.

Next, we test the procedure on a GAN trained on the CelebA dataset (Liu et al., 2015) at $64 \times 64$ resolution. Figure 3(a) shows the trajectories $\gamma_1$ starting at four fixed embeddings $z_1, \dots, z_4$. Although the association between the ordinal $k$ of the eigenvector $v_k$ and the attribute of the image being varied is not consistent throughout latent space, local changes still tend to occur along only one or two attributes at a time. Figure 3(b) shows the trajectories $\gamma_1, \gamma_2, \gamma_3$ and $\gamma_5$, all starting from the same fixed embedding $z_5$. As is the case for the dSprites dataset, we can see that trajectories $\gamma_k$ for distinct $k$ tend to effect changes to $G(z_0)$ along distinct attributes. These results suggest that, along trajectories from $z_0$ determined by a leading eigenvector of $M_z(z_0)$, changes in the output tend to occur along isolated attributes.

## 5. Aligning the Local Disentangled Factors

The $\beta$-VAE (Higgins et al., 2016) is known to learn disentangled representations, so that traveling along paths of the form
$$\gamma : t \mapsto z_0 + te_j, \tag{5}$$
for certain coordinate directions $e_j$, produces changes to isolated attributes of $G(z_0)$. In the previous section, we saw

---

**Algorithm 1** Procedure to trace path determined by $k$th leading eigenvector.

**Require:** $\mathrm{mv} : \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}^d$, $z, v \mapsto M_z(z)v$ is a function that computes matrix-vector products with the implicitly-defined matrix $M_z(z) \in \mathbb{R}^{m \times m}$.
**Require:** $z \in \mathbb{R}^m$ is the embedding from which to begin the trajectory.
**Require:** $k \in [1, m]$ is ordinal of the eigenvector to trace, with $k = 1$ corresponding to the leading eigenvector.
**Require:** $\alpha > 0$ is the step size.
**Require:** $\rho \in [0, 1)$ is the decay factor.
**Require:** $N \geqslant 1$ is the required number of steps in the trajectory.

1
2 **procedure** TRACEEIGENPATH($\mathrm{mv}, z, k, \alpha, \rho, N$)
3     $z_0 \leftarrow z$
4     **for** $i \in [1, N]$ **do**
5         $M_z(z) \leftarrow$ EVALUATENORMALJACOBIAN($\mathrm{mv}, z$)   ▷ Details in Appendix A
6         $V, D, V^t \leftarrow$ SVD($M_z(z)$)
7         $w_i \leftarrow v_k$                  ▷ Take the $k$th eigenvector
8         **if** $i \geqslant 2$ **then**
9             **if** $\langle w_{i-1}, w_i \rangle < 0$ **then**
10                 $w_i \leftarrow -w_i$   ▷ Prevent backtracking by ensuring that $\angle(w_{i-1}, w_i) \leqslant \pi/2$
11             $w_i \leftarrow \rho w_{i-1} + (1 - \rho) w_i$   ▷ Apply decay to smoothen the trajectory
12         $z_i \leftarrow z_{i-1} - \alpha w_i$
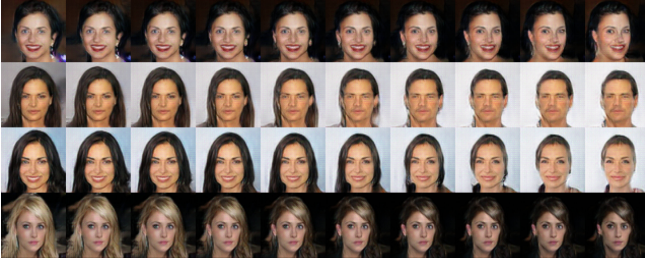13     **return** $\{z_0, \dots, z_N\}$

---

that such paths still exist for GANs, but they are not simply straight lines oriented along the coordinate axes (see Figure 2). We develop an efficient regularizer that encourages these paths to take the form of Equation 5, based on alignment of the top $k$ eigenvectors of $M_z(z_0)$ with the first $k$ coordinate directions $e_1, \dots, e_k$. Before proceeding, we describe a useful visualization technique to help measure the extent to which this happens. Let $M_z(z_0) = VDV^t$ be the eigendecomposition of $M_z(z_0)$, where $V$ is an orthogonal matrix whose columns are eigenvectors, and $D$ the diagonal matrix of nonnegative eigenvalues, sorted in descending order. Now we define $\tilde{V} : z \mapsto V$ to be the function that maps $z$ to the corresponding eigenvector matrix $V$, and let

$$F := \mathrm{E}_{z \sim p_z} \tilde{V}(z) \circ \tilde{V}(z), \tag{6}$$

where $p_z$ is the prior over $Z$, and '$\circ$' denotes Hadamard product. If the $k$th column is close to a one-hot vector, with values close to zero everywhere except at entry $j$, then we know that on average, the $k$th leading eigenvector of $M_z(z)$ is aligned with $e_j$. A heatmap generated from this matrix therefore allows us to gauge the extent to which each eigenvector $v_k$ is aligned with the coordinate direction $e_k$. Figure 5 shows that this does not happen automatically for a GAN, even when it is trained with a small latent variable size. Interestingly, eigenvector alignment does occur for $\beta$-VAEs. Appendix B explores this connection in more detail.

We begin by considering the case where we only seek to align the leading eigenvector $v_1 \in \mathbb{R}^m$ with the first coordinate direction $e_1 \in \mathbb{R}^m$. A simple way to do this is to obtain an estimate $\hat{v}_1$ for $v_1$ using $T$ power iterations, and then renormalize $\hat{v}_1 \in \mathbb{R}^m$ to a unit vector. We can then maximize the value of the first component of the elementwise squared vector $\hat{v}_1 \circ \hat{v}_1$, and minimize values of the remaining components. Using the mask $s_1 := (-1, 1, \dots, 1) \in \mathbb{R}^m$, we define the regular-

(a) Trajectories corresponding to the principal eigenvector for embeddings $z_1, \ldots, z_4$, from top to bottom, respectively. For embeddings $z_1$, $z_2$, and $z_3$, we show iterates $0, 20, \ldots, 180$, and for embedding $z_4$, we show iterates $0, 10, \ldots, 90$. The first trajectory varies azimuth; the second, gender; the third, hair length; and the fourth, hair color.



(b) Trajectories corresponding to the first, second, third, and fifth leading eigenvectors of embedding $z_5$, from top to bottom, respectively; we show iterates $0, 20, \ldots, 200$. The first trajectory primarily varies azimuth; the second, gender; the third, age; and the fourth, hair color and presence of facial hair.

Figure 3: Trajectories found by following leading eigenvectors from five fixed embeddings $z_1, \ldots, z_5$, using Algorithm 1 ($\alpha = 1.5 \cdot 10^{-2}$, $\rho = 0$). Details regarding the model architecture are given in Appendix C.

izer $R_1 : \mathbb{R}^m \to \mathbb{R}$,

$$R_1(z) := \sum_{i \in [1,m]} (s_1 \circ \hat{v}_1 \circ \hat{v}_1)_i. \tag{7}$$

Since $\hat{v}_1$ is constrained to unit norm, this regularizer is bounded. It can be incorporated into the loss for the generator using an appropriate penalty weight $\lambda > 0$.

Next, we consider the case where we would like to align the first two leading eigenvectors $v_1, v_2 \in \mathbb{R}^m$ with the first two coordinate directions $e_1, e_2 \in \mathbb{R}^m$. One potential approach is to first compute an estimate $\hat{v}_1$ to $v_1$ using $T$ power iterations, as before. Then, we could apply a modified version of the power method to obtain an estimate $\hat{v}_2$ for $v_2$, in which we project the result of each power iteration onto $\mathrm{span}(\hat{v}_1)^\perp$ using the projection $P_1 := I - \hat{v}_1 \hat{v}_1^t$. There are two problems with this approach. Firstly, it could be inaccurate: unless $\|v_1 - \hat{v}_1\| < \tau$ for sufficiently small $\tau > 0$, which may require a large number of power iterations, $P_1$ will not be an accurate projection onto $\mathrm{span}(v_1)^\perp$. Error in approximating $v_1$ would then jeopardize the approximation to $v_2$. Second, the approach is inefficient. We can only run the power method to estimate $v_2$ after the we have already obtained an estimate for $v_1$. If we use $T$ power iterations to estimate each eigenvector, then estimating the first $k$ eigenvectors will require a total of $kT$ power iterations. This is too slow to be practical.

**Algorithm 2** Procedure to estimate the top $k$ eigenpairs of $M_z(z)$.

**Require:** mv $: \mathbb{R}^m \to \mathbb{R}^m$ is a function that computes matrix-vector products with the implicitly-defined matrix $M_z(z) \in \mathbb{R}^{m \times m}$.
**Require:** $V \in \mathbb{R}^{m \times k}$ is a matrix whose columns are the initial estimates for the eigenvectors.
**Require:** $T \geqslant 1$ is the required number of power iterations.

```
 1
 2  ε ← 10⁻⁸                              ▷ Guards against division by numbers close to zero
 3
 4  procedure ESTIMATELEADINGEIGENPAIRS(mv, V, T)
 5      Let Mₖ ∈ ℝᵐˣᵏ be given by Equation 11
 6      V₀ ← Mₖ ∘ V                       ▷ '∘' denotes Hadamard product
 7
 8      for i ∈ [1, T] do
 9          Vᵢ ← M ∘ mv(Vᵢ₋₁)
10          Λᵢ ← diag(COLUMNNORMS(Vᵢ))
11          Vᵢ ← Vᵢ (Λᵢ + εI)⁻¹           ▷ Renormalize columns
12      return Λ_T, V_T
13
14  procedure COLUMNNORMS(A)                              ▷ A ∈ ℝᵖˣ�q
15      return (‖a₁‖, …, ‖a_q‖)           ▷ aᵢ ∈ ℝᵖ is the ith column of A
```

Our specific application of the power method enables an optimization that allows for the $k$ power iterations to be run simultaneously. Once we apply the regularizer to the GAN training procedure, the first eigenvector $v_1$ will quickly align with the first coordinate direction $e_1$. We therefore assume that $v_1 = e_1$. This assumption would imply that $\mathrm{span}(v_1)^\perp = \mathrm{span}(e_2, \ldots, e_m)$, so applying $P_1$ would amount to zeroing out the first component of $\hat{v}_2$ after each power iteration. Since $P_1$ would no longer depend on $\hat{v}_1$, we can run the power iterations for $v_1$ and $v_2$ in parallel. To formally describe this, we let $c_p := 1 \in \mathbb{R}^p$ be the constant vector of ones, and let

$$M_2 := \begin{pmatrix} 0 \\ c_m & c_{m-1} \end{pmatrix}. \tag{8}$$

Given a matrix $\hat{V}_t^{(2)} \in \mathbb{R}^{m \times 2}$ whose columns are the current estimates for $v_1$ and $v_2$, respectively, we can describe the power iterations for $v_1$ and $v_2$ using the recurrence

$$\hat{V}_{t+1}^{(2)} := M_2 \circ (M_z(z_0) \hat{V}_t^{(2)}). \tag{9}$$

Now, let $\hat{V}^{(2)} \in \mathbb{R}^{m \times 2}$ be the final estimate for $v_1$ and $v_2$. To implement the regularizer, we let $s_2 := (1, -1, 1, \ldots, 1) \in \mathbb{R}^m$, $S_2 := (s_1 \; s_2) \in \mathbb{R}^{m \times 2}$, and define

$$R_2(z) := \sum_{\substack{i \in [1,m] \\ j \in [1,2]}} (S_2 \circ \hat{V}^{(2)} \circ \hat{V}^{(2)})_{i,j}. \tag{10}$$

It is straightforward to generalize this approach to the case where we seek to align the first $k$ eigenvectors $v_1, \ldots, v_k$ with $e_1, \ldots, e_k$. For each eigenvector $v_i$, with $i \in [2, k]$, we assume that eigenvectors $v_1, \ldots v_{i-1}$ are already aligned with $e_1, \ldots, e_{i-1}$. The projections onto $\mathrm{span}(e_1)^\perp$, $\mathrm{span}(e_1, e_2)^\perp, \ldots, \mathrm{span}(e_1, \ldots, e_{i-1})^\perp$ can be implemented using columns $2, 3, \ldots, i$, respectively,

**Algorithm 3** Procedure to evaluate the alignment penalty.

**Require:** $k \in [1, m]$ is number of leading eigenvectors to align with $e_1, \ldots, e_k$.
**Require:** mv, $T$ are as defined in Algorithm 2.

```
1
2  procedure EVALUATEALIGNMENTREGULARIZER(k, mv, T)
3      Let S_k be given by Equation 12
4
5      α ← 2/(k(k + 1))
6      A ← diag(α · (k, k − 1, . . . , 1))
7      S_k ← S_k A          ▷ Reweight columns to prioritize alignment of leading
   eigenvectors
8
9      V_0 ← RANDOMRADEMACHER(m, k)
10     Λ̂, V̂ ← ESTIMATELEADINGEIGENPAIRS(mv, V_0, t)
11     return SUM(S_k ∘ V̂ ∘ V̂)
12
13 procedure RANDOMRADEMACHER(p, q)
14     return A ∈ ℝ^{p×q}, where a_{ij} = 1 with probability 1/2 and −1 with
   probability 1/2
```
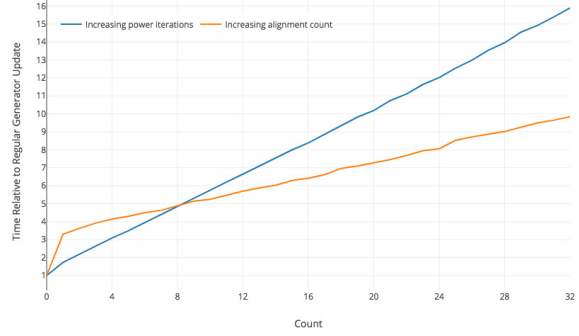


Figure 4: Median relative cost per generator update as a function of the number of eigenvectors $k$ to align, and the number of power iterations $T$ to use for Algorithm 3. The cost is measured relative to the time required to make one RMSProp (Tieleman & Hinton, 2012) update to the parameters of a DCGAN generator (approximately $20.419$ ms) with base feature map count 64 and batch size 32. Complete details regarding the model architecture are given in Appendix C. Medians were computed using the update times for the first 500 iterations; the median absolute deviations are small enough as to be indistinguishable from the medians on the plot.
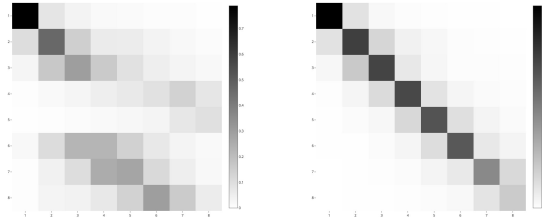
of the mask $M_k \in \mathbb{R}^{m \times k}$. This mask, which is a generalization of the one defined by Equation 8, is given by $\mathbb{R}^{m \times k} \ni M_k :=$

$$\begin{pmatrix} c_k & c_k - e_1 & c_k - (e_1 + e_2) & \cdots & c_k - (e_1 + \cdots + e_{k-1}) \\ c_{m-k} & c_{m-k} & c_{m-k} & \cdots & c_{m-k} \end{pmatrix}. \quad (11)$$

The resulting procedure to estimate the leading $k$ eigenvectors is described by Algorithm 2. Figure 4 shows that the runtime of Algorithm 2 scales linearly with respect to the number of eigenvectors $k$ and the number of power iterations $T$. Next, we generalize Equation 10, in order to describe how to evaluate the regularizer $R_k$. Let $s_p \in \mathbb{R}^m$ be given by $(s_p)_i = -1$ if $i = p$ and 1 otherwise, and define

$$S_k := \begin{pmatrix} -1 & 1 & 1 & \cdots & 1 \\ 1 & -1 & 1 & \cdots & 1 \\ 1 & 1 & -1 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & -1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{pmatrix} = \begin{pmatrix} s_1 \cdots s_k \end{pmatrix} \in \mathbb{R}^{m \times k}. \quad (12)$$

Algorithm 3 shows how $S_k$ is used with the result of Algorithm 2 to evaluate $R_k$. Finally, Figure 2(c) shows that incorporating the alignment regularizer into the generator loss successfully aligns the trajectories produced by Algorithm 1 with the coordinate axes.

If the alignment regularizer is implemented exactly as described by Equation 10, it will fail to have the intended effect. The reason for this has to do with the assumption behind the optimization used to run the $k$ power iterations in parallel. Before attempting to align $v_i$ with $e_i$, we assume that that $v_1, \ldots, v_{i-1}$ are already aligned with $e_1, \ldots, e_{i-1}$. When this assumption fails to hold, the projections computed using the columns of $M_k$ will no longer be valid. Figure 5(a) shows that the matrix $F$ does not have a diagonal in its top-left corner, which is what we would expect to see if $v_1, \ldots, v_k$ were aligned with $e_1, \ldots, e_k$. Fortunately, there is a simple fix that remedies the situation. We would like to encourage the optimizer to prioritize alignment of $v_i$ with $e_i$ over alignment of $v_{i+1}, \ldots, v_k$ with $e_{i+1}, \ldots, e_k$,



Figure 5: Comparison of the top-left $8 \times 8$ corner of the matrix $F$ (Equation 6) for alignment-regularized GANs ($k = 8$, $T = 8$, $\lambda = 0.1$), trained with reweighting (left) and without reweighting (right) of the columns of the matrix $S_k$ defined by Equation 12. Both GANs were trained on the CelebA dataset (Liu et al., 2015) at $64 \times 64$ resolution, with latent variable size 128. See Appendix C for complete details regarding the model architecture and training procedure.

for all $i \in [1, k - 1]$. A simple way to do this is to multiply the $i$th column of $M_k$ by a weight of $(k - i + 1)\alpha$. We choose $\alpha$ based on the condition that these weights sum to one, i.e.,

$$\sum_{i \in [1, k]} i\alpha = 1, \quad \text{implying that} \quad \alpha = \frac{2}{k(k+1)}. \quad (13)$$

This reweighting scheme is implemented in lines 5–7 of Algorithm 3. Figure 5(b) confirms that this modification induces the desired structure in the top-left corner of $F$.

## 6. Results

We first make a quantitative comparison between our approach and previous methods that have been used to obtain disentangled representations. This requires us to measure the extent to which the second property of disentangled representations – namely, interpretability of changes resulting from perturbations to individual coordinates of a latent variable – holds. Suppose that we had knowledge of the ground truth latent factors for the dataset, along with a simulator

| Model | Disentanglement Score |
|---|---|
| Ground truth | 100 |
| Raw pixels | $45.75 \pm 0.8$ |
| PCA | $84.90 \pm 0.4$ |
| ICA | $42.03 \pm 10.6$ |
| DC-IGN | $\mathbf{99.3} \pm \mathbf{0.1}$ |
| InfoGAN[1] | $73.5 \pm 0.9$ |
| VAE (untrained) | $44.14 \pm 2.5$ |
| VAE | $61.58 \pm 0.5$ |
| $\beta$-VAE | $\mathbf{99.23} \pm \mathbf{0.1}$ |
| Ours | $\mathbf{92.34} \pm \mathbf{0.4}$ |

Table 1: Comparison of disentanglement metric scores for our method to those reported in (Higgins et al., 2016). We use an alignment-regularized GAN ($k, \lambda = 6, 0.6$) with latent variable size 10. Details regarding the model architecture and training procedure are given in Appendix C.

that can synthesize new outputs given assignments to these latent factors. Then we could generate pairs of outputs, such that the outputs in each pair differ only along a single attribute. Let $(x_0, x_1)$ be one such pair, and $z_0 := G^{-1}(x_0)$ and $z_1 := G^{-1}(x_1)$ the corresponding latent variables for the generator $G : \mathbb{R}^m \to \mathbb{R}^n$. If $G$ satisfies the second property, then we would expect $z_0$ and $z_1$ to be approximately equal along all components, except the one corresponding to the attribute that was varied. Hence, $|z_0 - z_1|$ should be a one-hot vector in expectation. (Higgins et al., 2016) propose an evaluation metric based on this idea. It involves training a linear classifier to predict the index of the latent factor that was varied in order to generate each pair. At each step of training for the classifier, we sample a batch of input-target pairs in accordance with the procedure described in (Higgins et al., 2016), and update the classifier using the cross-entropy loss.

Application of the evaluation metric to GANs is complicated by the fact that a direct procedure to invert the generator is usually not available. Models such as the $\beta$-VAE consist of an encoder that effectively functions as an inverse for the generative model, so this is not an issue. For the purpose of evaluation, we also train an encoder to invert the fixed generator, after GAN training has finished. This additional training procedure uses a standard autoencoding loss, and the details are specified in Appendix D. Table 1 compares our approach to the ones evaluated in (Higgins et al., 2016) on the dSprites dataset (Matthey et al., 2017), which we describe in Section 4. As stated in Appendix C, we added noise to both the real and fake inputs of the discriminator in order to stabilize GAN training on this dataset, for which the pixels are binary-valued. Incorporating this modification into the training procedure for InfoGAN may improve the score reported by (Higgins et al., 2016). Our approach achieves a score competitive to that of DC-IGN, which makes use

of supervised information. We note that the encoder trained to invert the generator does not always succeed in finding a latent variable that yields an accurate reconstruction of the input. This may lead to a reduction in the disentanglement score, a problem that does not affect the other approaches for which an accurate inference procedure is available.

Next, we make a qualitative comparison between our approach and $\beta$-VAE. We train a series of GANs on the CelebA dataset (Liu et al., 2015), with $k \in \{8, 16, 32\}$ for Algorithm 3 and varying values for the alignment regularizer weight $\lambda$. We also train a series of $\beta$-VAEs with $\beta \in \{1, 2, 4, 8, 16\}$. The results for our method with $k, \lambda = 32, 0.5$ are shown in Figure 7, and the results for $\beta$-VAE with $\beta = 8$ in Figure 8. Figure 6 shows the top-left corners of the matrix $F$ given by Equation 6, for three configurations of our approach with $k \in \{8, 16, 32\}$. A penalty weight of $\lambda = 0.5$ was not sufficient to result in a clear diagonal structure in the top-left corner of the matrix $F$ for the configuration with $k = 32$. Nonetheless, this configuration resulted in the largest number of disentangled factors, and the results are shown in Figure 7. The best results for the configurations with $k \in \{8, 16\}$ are shown in Appendix E. We found that the $\beta$-VAE configuration with $\beta = 8$ resulted in the largest number of disentangled factors, while still maintaining sufficiently low reconstruction error so as to keep the sample quality acceptable. These results are shown in Figure 8. In addition to better sample quality, our approach is able to learn concepts such as different kinds of hair styles (coordinates 9 and 10 in Figure 7) that are not modeled by the $\beta$-VAE.

Even with relatively large values for the penalty weight, the heatmaps shown in Figure 6 contain nonzero entries above and below the diagonal. This can result in some degree of leakage of the attribute controlled by one coordinate into the next. To see this in more detail, we examine the disentangled factors found by the configuration with $k = 16$, which are shown in Figure 12. Coordinates $\{1, 2, 3\}$ all involve change in hair darkness; coordinates $\{4, 5, 6\}$ all involve change in gender; coordinates $\{6, 7, 8, 9\}$ all involve smiling; and coordinates $\{12, 13\}$ all involve change in the location of the hair partition. This is a limitation of our current approach, and modifications to the implementation of the alignment regularizer in Algorithm 3 may help to mitigate this leakage. We plan to investigate such improvements in future work.

## 7. Conclusion

Our work approaches the problem of learning disentangled representations from the perspective of eigenvector alignment. We develop a novel regularizer which, when incorporated into the GAN objective, induces disentangled representations of quality comparable to those obtained by VAE-based approaches (Higgins et al. (2016), Kim & Mnih
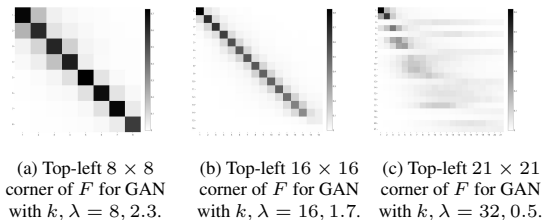
(a) Top-left $8 \times 8$ corner of $F$ for GAN with $k, \lambda = 8, 2.3$.

(b) Top-left $16 \times 16$ corner of $F$ for GAN with $k, \lambda = 16, 1.7$.

(c) Top-left $21 \times 21$ corner of $F$ for GAN with $k, \lambda = 32, 0.5$.

Figure 6: Comparison of the top-left corners of the matrix $F$ (Equation 6) for alignment-regularized GANs with $k \in \{8, 16, 32\}$. All GANs were trained with latent variable size 128. Details regarding the model architecture and training procedure are given in Appendix C.

(2018), Chen et al. (2018), Esmaeili et al. (2018)), without the need to introduce auxiliary models into the training procedure (Chen et al., 2016). This approach is also not specific to the GAN framework (Goodfellow et al., 2014), and could potentially be applied to autoregressive models, such as those used to generate text and audio. We believe this is an important direction for future investigation. So far, two different perspectives for viewing disentanglement have been proposed: maximizing mutual information and eigenvector alignment. An investigation into the relationship between the two could further our understanding of what disentanglement is and why it occurs.

## Acknowledgements

| Coordinate | Description | Example |
|---|---|---|
| 1 | Background darkness | |
| 2 | Azimuth | |
| 3 | Bangs | |
| 4 | Gender, bangs, and smiling | |
| 5 | Smiling | |
| 6 | Hair color | |
| 7 | Hair color and hair style | |
| 8 | Lighting color and bangs | |
| 9 | Hair color and hair style | |
| 10 | Hair color and hair style | |
| 11 | Jawline | |
| 12 | Smiling and bangs | |
| 13 | Background color | |
| 14 | Age and hairline | |
| 15 | Age | |
| 17 | Location of hair partition | |
| 18 | Lighting and skin tone | |
| 20 | Mouth open | |
| 21 | Mouth open | |



Figure 7: Disentanglement results for alignment-regularized GAN ($k, \lambda = 32, 0.5$) with latent variable size 128 on the CelebA dataset, at $64 \times 64$ resolution. Details regarding the model architecture and training procedure are given in Appendix C. Results with additional samples for each coordinate can be found at this URL: https://drive.google.com/open?id=1E2TneLSAdyFYgN4GUzYKnHDYABo-G8RK.

| Coordinate | Description | Example |
|---|---|---|
| 1 | Azimuth | |
| 2 | Background color | |
| 3 | Location of hair partition | |
| 4 | Skin color | |
| 5 | Hair direction | |
| 8 | Background color | |
| 12 | Hair length | |
| 14 | Jawline | |
| 19 | Skin tone | |
| 20 | Hair style | |
| 22 | Sunglasses | |
| 23 | Gender | |
| 25 | Background darkness | |
| 26 | Hairline and skin tone | |
| 28 | Lighting direction and rotation | |
| 32 | Hair size, color, and smiling | |



Figure 8: Disentanglement results for $\beta$-VAE ($\beta = 8$) with latent variable size 32 on the CelebA dataset, at $64 \times 64$ resolution. Details regarding the model architecture and training procedure are given in Appendix C. Results with additional samples for each coordinate can be found at this URL: https://drive.google.com/open?id=1Zu5H_M0dOuE2NYeOJ9WZFSux4U5qM5Pq.

# References

Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pp. 265–283, 2016.

Aubry, M., Maturana, D., Efros, A. A., Russell, B. C., and Sivic, J. Seeing 3d chairs: exemplar part-based 2d-3d alignment using a large dataset of cad models. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3762–3769, 2014.

Baydin, A. G., Pearlmutter, B. A., Radul, A. A., and Siskind, J. M. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*, 18(153): 1–153, 2017.

Brock, A., Donahue, J., and Simonyan, K. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.

Burgess, C. P., Higgins, I., Pal, A., Matthey, L., Watters, N., Desjardins, G., and Lerchner, A. Understanding disentangling in $\beta$-vae. *arXiv preprint arXiv:1804.03599*, 2018.

Chen, T. Q., Li, X., Grosse, R., and Duvenaud, D. Isolating sources of disentanglement in variational autoencoders. *arXiv preprint arXiv:1802.04942*, 2018.

Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., and Abbeel, P. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, pp. 2172–2180, 2016.

Esmaeili, B., Wu, H., Jain, S., Bozkurt, A., Siddharth, N., Paige, B., Brooks, D. H., Dy, J., and van de Meent, J.-W. Structured disentangled representations. *stat*, 1050:29, 2018.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.

Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016.

Karras, T., Aila, T., Laine, S., and Lehtinen, J. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.

Kim, H. and Mnih, A. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.

Matthey, L., Higgins, I., Hassabis, D., and Lerchner, A. dsprites: Disentanglement testing sprites dataset. https://github.com/deepmind/dsprites-dataset/, 2017.

Mescheder, L., Geiger, A., and Nowozin, S. Which training methods for gans do actually converge? In *International Conference on Machine Learning*, pp. 3478–3487, 2018.

Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.

Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.

Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.

Salimans, T. and Kingma, D. P. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 901–909, 2016.

Tieleman, T. and Hinton, G. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.

Townsend, J. A new trick for calculating jacobian vector products. `https://j-towns.github.io/2017/06/12/A-new-trick.html`, 2017.

Xiang, S. and Li, H. On the effects of batch and weight normalization in generative adversarial networks. *stat*, 1050:22, 2017.

## A. Forward- and reverse-mode automatic differentiation

The entries of $J_G$ are not explicitly stored in memory: it is an implicit-defined matrix. As such, information about $J_G$ must be accessed via automatic differentiation (AD), of which there are two kinds: forward-mode and reverse-mode. We briefly describe them here, and refer the reader to (Baydin et al., 2017) for a more comprehensive survey.

Suppose that we are given a differentiable function $f : \mathbb{R}^m \to \mathbb{R}^n$ corresponding to a feedforward neural network with $L$ layers, and let $G := (V, E)$ be its representation as a computation graph. For each $k \in [1, L]$, let $R_k$ be the set of nodes corresponding to the $k$th layer, and let $r_k : \mathbb{R}^{\alpha_{k-1}} \to \mathbb{R}^{\alpha_k}$ be the function computed by the nodes in $R_k$, with $\alpha_0 := m$ and $\alpha_L := n$. Then $f = r_L \circ r_{L-1} \circ \cdots \circ r_2 \circ r_1$. Finally, let $\bar{x} \in \mathbb{R}^m$ denote the variable for the input to $f$, and $\bar{b}_{k-1} \in \mathbb{R}^{\alpha_{k-1}}$ the variable for the input to $r_k$, with $\bar{b}_0 := \bar{x}$. Given a vector $v \in \mathbb{R}^m$, forward-mode AD computes $Jv$. It works according to the recurrence

$$v_k := \left( D_{\bar{x}} r_k \big|_{\bar{x}=x} \right) v = \left( D_{\bar{b}_{k-1}} r_k \big|_{\bar{b}_{k-1}=b_{k-1}} \right) \left( D_{\bar{x}}\, r_{k-1} \circ \ldots \circ r_1 \big|_{\bar{x}=x} \right) v$$
$$= \left( D_{\bar{b}_{k-1}} r_k \big|_{\bar{b}_{k-1}=b_{k-1}} \right) v_{k-1},$$

where $k \in [1, L]$ and $v_0 := v$. After step $k$, we will have obtained the product of the Jacobian of $r_k$ with respect to $x$, and $v$; after step $L$, we will have obtained the desired product $Jv$.

Forward-mode AD computes $Jv$ far more efficiently than the approach of first evaluating $J$, and subsequently multiplying by $v$. Suppose for simplicity that $\alpha_k = n$ for all $k \in [0, L]$, so that $J \in \mathbb{R}^{n \times n}$. To compute $J$ independently, we would use the chain rule, which gives

$$D_{\bar{x}} f \big|_{\bar{x}=x} = \left( D_{\bar{b}_{L-1}} r_L \big|_{\bar{b}_{L-1}=b_{L-1}} \right) \left( D_{\bar{b}_{L-2}} r_{L-1} \big|_{\bar{b}_{L-2}=b_{p-2}} \right) \cdots \left( D_{\bar{b}_0} r_1 \big|_{\bar{b}_0=x} \right).$$

Each layer $r_k$, for $k \in [1, L]$, must compute its $n \times n$ Jacobian, and each layer after the first must multiply its Jacobian with the $n \times n$ Jacobian of the preceding layer with respect to $x$. This process is described by the recurrence

$$D_{\bar{x}} r_k \big|_{\bar{x}=x} = \left( D_{\bar{b}_{k-1}} r_k \big|_{\bar{b}_{k-1}=b_{k-1}} \right) \left( D_{\bar{x}} r_{k-1} \big|_{\bar{x}=x} \right).$$

Assuming that $\Theta(1)$ operations are required to compute each element of the Jacobian, the total number of operations required is proportional to

$$N := n^2 + (L-1)(n^2 + n^3) \in \Theta(n^3),$$

despite the fact that $J$ has only $n^2$ elements. On the other hand, the $k$th step of the forward-mode recurrence requires $\Theta(n)$ operations, so computing $Jv$ only requires $\Theta(Ln)$ operations. By running forward-mode with $v = e_i$ for each $i \in [1, n]$, we can form $J$ column-by-column in only $\Theta(Ln^2)$ operations. Henceforth, we will measure operation count in terms of invocations to the AD engine, rather than by counting elementary operations.

Of the two kinds of AD, it is reverse-mode that finds the most use in machine learning. Given a vector $w \in \mathbb{R}^n$, it computes $J^t w$. Most applications involve minimizing a scalar-valued loss function with respect to a vector of parameters, which corresponds to the case where $n = 1$. This special case is otherwise known as backpropagation. Unlike forward-mode, which begins at the first layer and ends at the last, reverse-mode begins at the last layer and ends at the first. It works according to the recurrence

$$w_k^t := w^t \left( D_{\bar{b}_{k-1}} r_k \big|_{\bar{b}_{k-1}=b_{k-1}} \right) = w^t \left( D_{\bar{b}_k}\, r_L \circ \cdots \circ r_{k+1} \big|_{\bar{b}_k=b_k} \right) \left( D_{\bar{b}_{k-1}} r_k \big|_{\bar{b}_{k-1}=b_{k-1}} \right)$$
$$= w_{k+1}^t \left( D_{\bar{b}_{k-1}} r_k \big|_{\bar{b}_{k-1}=b_{k-1}} \right),$$

where $k$ ranges from $L$ to 1, and $w_L := w$. After $L$ steps, we will have obtained the desired product $J^t w$. Machine learning frameworks typically expose the full interface to the reverse-mode AD engine, rather than specializing to the case of backpropagation. E.g., in TensorFlow (Abadi et al., 2016), one can change the value of $w$ for `tf.gradients` from a vector of ones to a custom value specified by the `grad_ys` parameter.

If desired, we can compute the entire Jacobian using either type of AD. By running forward-mode with $v = e_i$ for each $i \in [1, m]$, we can form $J$ column-by-column, using $m$ total invocations to AD. Similarly, by running reverse-mode

```
1    import tensorflow as tf
2
3    def forward_gradients(ys, xs, d_xs):
4        v = tf.placeholder_with_default(tf.ones_like(ys), shape=ys.get_shape())
5        g = tf.gradients(ys, xs, grad_ys=v)
6        return tf.gradients(g, v, grad_ys=d_xs)
7
8    def j_v(ys, xs, vs):
9        return forward_gradients(ys, xs, vs)
10
11   def jt_v(ys, xs, vs):
12       return tf.gradients(ys, xs, vs)
13
14   def jt_j_v(ys, xs, vs):
15       jv = j_v(ys, xs, vs)
16       return tf.gradients(ys, xs, jv)
17
18   def j_jt_v(ys, xs, vs):
19       jt_v_ = jt_v(ys, xs, vs)
20       return forward_gradients(ys, xs, jt_v_)
```

Listing 1: TensorFlow implementation of Jacobian-vector operations.



(a)                                      (b)

Figure 9: Comparison between the average squared eigenvector matrices $F$ (Equation 6) for a GAN generator with latent variable size 16 (left) and the decoder from a $\beta$-VAE ($\beta = 8$) with latent variable size 32 (right), both trained on the CelebA dataset. Complete architecture and training details are provided in Appendix C. The matrix $F$ for the VAE decoder contains several columns that are close to one-hot vectors. By contrast, the same matrix for the GAN generator exhibits little structure, despite the fact that it is trained with a small latent variable size.

with $w = e_i$ for each $i \in [1, n]$, we can form $J$ row-by-row using $n$ total invocations to AD. If $n > m$, the former is typically faster; otherwise, the latter is preferable. We note that while AD offers a relatively efficient approach for evaluating $J$, doing so at each iteration of optimization becomes impractical.

Many popular machine learning frameworks (e.g., TensorFlow) do not implement forward-mode natively, since it is seldom used for machine learning. Surprisingly, this is not a limitation: *reverse-mode can be used to implement forward-mode.* Given a differentiable function $f : \mathbb{R}^m \to \mathbb{R}^n$, we can compute $w^t J$ for a given vector $w \in \mathbb{R}^n$, using reverse-mode AD. Treating the input to $f$ as a constant, we can regard $w^t J$ as a function $g : \mathbb{R}^n \to \mathbb{R}^m$, $w \mapsto w^t J$. The derivative of $g$ with respect to $w$ is given by $J^t$, and so another application of reverse-mode AD allows us to compute $v^t J^t = (Jv)^t$. Hence, reverse-mode AD can be used to implement forward-mode AD. This trick was first described by (Townsend, 2017). We provide a TensorFlow implementation of the procedures to compute $Jv$, $J^t v$, $J^t Jv$, and $JJ^t v$ in Listing 1.

## B. Eigenvector Alignment for $\beta$-VAEs

The results from Section 4 suggest that alignment of the eigenvectors of $M_z(z_0)$ with the coordinate axes might be sufficient to induce disentanglement in the latent representation of a generative model. The $\beta$-VAE is known to learn such a representation when $\beta$ is increased, so that the KL-divergence between the approximate posterior and the prior of the decoder is made sufficiently small. Suppose that a $\beta$-VAE exhibits disentanglement along the $j$th coordinate direction in

Figure 10: Investigation of whether the matrix $F$ can be used to determine which directions in latent space correspond to disentangled factors for a $\beta$-VAE. For each $i \in [1, 32]$, we compute the maximum value along the $i$th row of matrix $F$ shown in Figure 9(b). Then, we plot a point indicating whether or not disentanglement occurs along this direction, as determined by visual inspection. The directions that do result in disentanglement are shown in Figure 8.

| Figure | DCGAN Base Feature Map Count | Latent Variable Size | Notes |
|--------|------------------------------|----------------------|-------|
| 13(c) | 64 | 32 | 1 |
| 13(e) | 64 | 64 | 1 |
| 13(g) | 64 | 128 | 1 |
| 13(d),5, 4, 6(a), 6(b), 7 | 64 | 128 | 1 |
| 3, 11, 12 | 128 | 128 | 1 |
| 13(f) | 64 | 256 | 1 |
| 13(h) | 64 | 512 | 1 |
| 2 | 64 | 3 | 1, 3 |
| 1 | 64 | 10 | 1, 3 |
| 9(a) | 64 | 16 | 1 |
| 9(b), 10 | 64 | 32 | 2 |

Table 2: GAN and $\beta$-VAE architectures used for all figures and tables. Note 1: weight normalization (Salimans & Kingma, 2016) was applied both to the generator and the discriminator, with the scale $g$ fixed to one for the discriminator. Note 2: spectral weight normalization (Miyato et al., 2018) was applied both to the encoder and the decoder, with learned scales for both. Note 3: gaussian noise with standard deviation 0.6 was added to both real and fake inputs to the discriminator.

latent space. Then paths of the form $\gamma : t \mapsto z_0 + te_j$ will produce changes to $G(z_0)$ along isolated factors of variation. If such a path coincides with the trajectory found by Algorithm 1 for the $k$th leading eigenvector throughout latent space, then this eigenvector must be aligned with the $j$th coordinate axis. Hence, we would expect the $k$th column of the matrix $F$ given by Equation 6 to be a one-hot vector close to $e_j$. Figure 9 shows that, in fact, several columns of $F$ have high similarity to coordinate directions.

Next, we investigate whether a column of $F$ having high similarity with a coordinate direction $e_j$ implies that disentanglement actually occurs along $\gamma : t \mapsto z_0 + te_j$. Figure 10 shows that with the exception of three coordinates having similarity greater than 0.35, this turns out not to be the case. In other words, several of the eigenvectors naturally align with coordinate directions during training, but disentanglement does not reliably occur along all of them. More work needs to be done in order to better understand the relationship between eigenvector alignment and disentangled representations.

## C. Generator Architectures and Training Procedure

All models in this work are based on the DCGAN (Radford et al., 2015) architecture. Table 2 describes the architectures of the generator and discriminator (in the case of GANs) and the encoder and decoder (in the case of VAEs) associated with each figure and table. All models use the translated PReLU activation function (Xiang & Li, 2017); the ReLU leaks are learned, and clipped to the interval $[0, 1]$ after each parameter update. The only data preprocessing we applied was to scale the pixel values to the interval $[0, 1]$. The GANs were trained using the original, non-saturating GAN loss described in (Goodfellow et al., 2014) with multivariate normal prior. Both the generator and the discriminator were trained using RMSProp (Tieleman & Hinton, 2012) with step size $10^{-4}$, decay factor $0.9$, and $\epsilon = 10^{-6}$. Each parameter update was made using a batch size of 32, and the models were trained for a total of 750,000 updates. The VAEs were trained using a gaussian likelihood model for the decoder, and the log-diagonal covariance parameterization for the encoder. The fixed, per-pixel standard deviation of the decoder was chosen such that, disregarding constant terms, the log-likelihood corresponds to the reconstruction error, normalized by the product of the number of channels and the number of pixels. To optimize the evidence lower bound, we used Adam (Kingma & Ba, 2014) with $\beta_1, \beta_2, \epsilon = 0.5, 0.99, 10^{-8}$. Each parameter update was made using a batch size of 32, and the models were trained for a total of 1,500,000 updates. The KL-divergence weight $\beta$ was increased linearly

---

**Algorithm 4** Procedure to generate a batch for disentanglement metric classifier.

---

**Require:** $n_{inst} \geqslant 1$ is the number of aligned pairs to use, in order to create each instance in the batch.
**Require:** $n_{batch} \geqslant 1$ is the batch size.
**Require:** $\text{Enc} : \mathbb{R}^n \to \mathbb{R}^m$ is the encoder functioning as the inverse of the generator $G$.

$c \leftarrow (3, 6, 40, 30, 30)$                                       ▷ Ranges for the five attributes that determine each shape

**procedure** SAMPLESHAPE($i, v$)                                      ▷ Sample shape with attribute $i$ fixed to value $v$
    **for** $j \in [1, 5]$ **do**
        $u_j \leftarrow \text{RANDOMUNIFORM}(1, c_j)$
    $u \leftarrow (u_1, u_2, u_3, u_4, u_5)$
    $u_i \leftarrow v$
    **return** MAKESHAPE($u$)

**procedure** MAKEINSTANCE($n_{inst}, \text{Enc}$)
    $i \leftarrow \text{RANDOMUNIFORM}(2, 5)$                    ▷ Sample index of attribute to fix
    $v \leftarrow \text{RANDOMUNIFORM}(1, c_i)$                   ▷ Sample value for fixed attribute
    $z \leftarrow 0 \in \mathbb{R}^m$
    **for** $k \in [1, n_{inst}]$ **do**
        $z_0 \leftarrow \text{Enc}(\text{SAMPLESHAPE}(i, v))$
        $z_1 \leftarrow \text{Enc}(\text{SAMPLESHAPE}(i, v))$
        $z \leftarrow z + |z_0 - z_1|$
    **return** $z/n_{inst}, i$

**procedure** MAKEBATCH($n_{inst}, n_{batch}, \text{Enc}$)
    inputs $\leftarrow \varnothing$
    targets $\leftarrow \varnothing$
    **for** $k \in [1, n_{batch}]$ **do**
        $x, y \leftarrow \text{MAKEINSTANCE}(n_{inst}, \text{Enc})$
        inputs $\leftarrow$ inputs $\cup \{x\}$
        targets $\leftarrow$ targets $\cup \{y\}$
    **return** inputs, targets

---

from 0 to the final value over the first 1,000,000 updates. We applied the function $x \mapsto 2(\tanh(x) + 1/2)$ to the outputs of both the GAN generators and the VAE decoders; this ensures that the output pixel values are in the interval $[0, 1]$. To predict the mean and log-diagonal covariance with the VAE encoders, we apply the translated ReLU activation function (Xiang & Li, 2017) to the final convolutional features, followed by two fully-connected layers, one for each statistic.

## D. Implementation Details for Disentanglement Metric Score

Algorithm 4 describes the procedure used to generate the batches to train and evaluate the classifier for the disentanglement metric (Higgins et al., 2016). We use $n_{inst}, n_{batch} = 64, 32$, and update the classifier using SGD with nesterov accelerated gradient (step size $10^{-2}$, momentum 0.99). We train the classifier using a total of 10,000 parameter updates, and evaluate its performance using 5000 instances, as reported by (Higgins et al., 2016). To invert the pretrained GAN generator, we train a VAE decoder with twice the base feature map count for the generator, using 30,000 parameter updates. The details for this training procedure are identical to those for regular VAE training described in Appendix C, except that the generator is not updated, and the KL-divergence weight $\beta$ is set to zero. In other words, we use a standard autoencoding loss with a fixed decoder.

# E. Additional GAN Disentanglement Results

| Coordinate | Description | Example |
|:---:|:---|:---|
| 1 | Background darkness and hair color |  |
| 2 | Azimuth, lighting, and hair color |  |
| 3 | Hairline and hair color |  |
| 4 | Azimuth |  |
| 5 | Shadow |  |
| 6 | Smiling, age, skin tone, gender |  |
| 7 | Smiling, age, jawline |  |
| 8 | Jawline and hairstyle |  |

Figure 11: Disentanglement results for alignment-regularized GAN ($k, \lambda = 8, 2.3$). Details regarding the model architecture and training procedure are given in Appendix C. Results with additional samples for each coordinate can be found at this URL: `https://drive.google.com/open?id=1Wnd5jIxopFsBRlylMUN2HfWhICXGiU3u`.

| Coordinate | Description | Example |
|:---:|:---|:---|
| 1 | Background and hair darkness | |
| 2 | Azimuth and hair darkness | |
| 3 | Hair length, hair darkness, and lighting | |
| 4 | Smiling and gender | |
| 5 | Hair length, hair darkness, and gender | |
| 6 | Smiling, bangs, and gender | |
| 7 | Smiling and hairstyle | |
| 8 | Smiling, jawline, and glaring expression | |
| 9 | Smiling, bangs, and mouth open | |
| 10 | Hairline | |
| 11 | Raised eyebrows and skin tone | |
| 12 | Raised eyebrows and location of hair partition | |
| 13 | Raised eyebrows and location of hair partition | |
| 14 | Lighting color | |

Figure 12: Disentanglement results for alignment-regularized GAN ($k, \lambda = 16, 1.7$). Details regarding the model architecture and training procedure are given in Appendix C. Results with additional samples for each corodinate can be found at this URL: `https://drive.google.com/open?id=1BM-P-hMF7sV_0smNFD_iTAKpq6FeNsCo`.
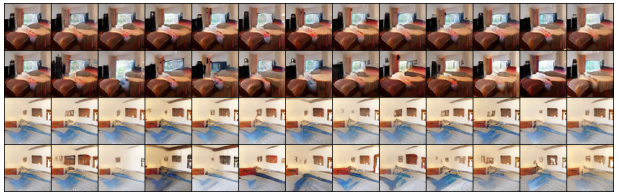
# F. Supplementary Figures

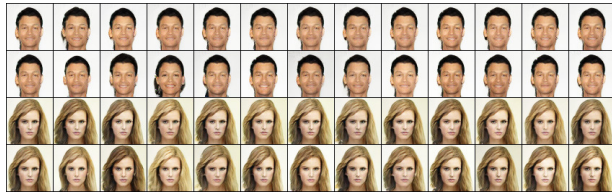(a) Log spectra for CelebA models with $n_f = 64$ and $m$ varied.



(b) Log spectra for LSUN Bedroom models with $n_f = 256$ and $m$ varied.
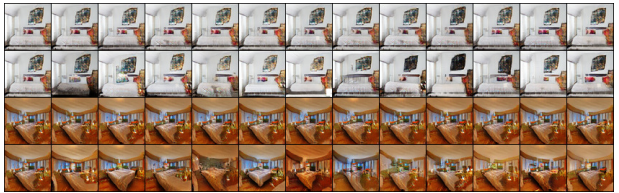


(c) CelebA ($m = 32, n_f = 64, \epsilon = 0.40$), embeddings $z_{10}, z_{11}$.
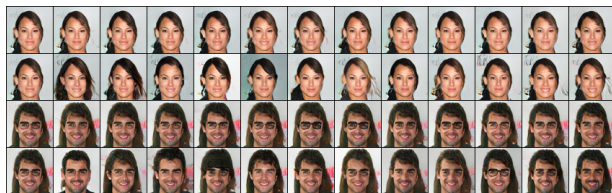


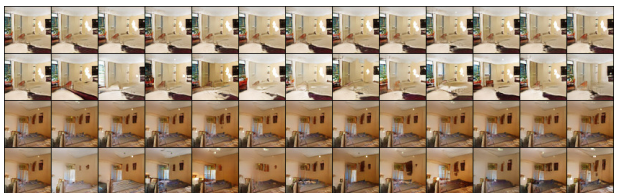(d) LSUN Bedroom ($m = 128, n_f = 64, \epsilon = 0.80$), embeddings $z_5, z_{12}$.



(e) CelebA ($m = 64, n_f = 64, \epsilon = 0.65$), embeddings $z_6, z_{12}$.



(f) LSUN Bedroom ($m = 256, n_f = 64, \epsilon = 0.80$), embeddings $z_6, z_8$.



(g) CelebA ($m = 128, n_f = 64, \epsilon = 0.80$), embeddings $z_7, z_9$.



(h) LSUN Bedroom ($m = 512, n_f = 64, \epsilon = 0.80$), embeddings $z_8, z_{10}$.

Figure 13: Effect of perturbing a latent variable $z$ along the leading eigenvectors of $M_z(z)$. Subfigures (a) and (b) show the top eigenvalues of $M_z(z_i)$, where $\{z_1, \ldots, z_{12}\}$ are fixed embeddings; small eigenvalues are omitted. Subfigures (c)–(d) compare the effects of perturbations along random directions to perturbations along leading eigenvectors. Each subfigure consists of two stacked two-row grids. The leftmost images of each grid are both identical and equal to $G(z_i)$, for some $i \in [1, 12]$. The first row shows the effect of perturbing $z_i$ along 13 directions sampled uniformly from the sphere of radius $\epsilon$, while the second row shows the effect of perturbing $z_i$ along the first 13 leading eigenvectors of $M_z(z_i)$, by the same distance $\epsilon$. Details regarding the model architectures are given in Appendix C.