

Disentangled Representation Learning and Generation With Manifold Optimization

Arun Pandey

arun.pandey@esat.kuleuven.be

KU Leuven, Department of Electrical Engineering, STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics, B-3001 Leuven, Belgium

Michaël Fanuel

michael.fanuel@univ-lille.fr

Université de Lille, CNRS, Centrale Lille, F-59000 Lille, France

Joachim Schreurs

joachim.schreurs@esat.kuleuven.be

Johan A. K. Suykens

johan.suykens@esat.kuleuven.be

KU Leuven, Department of Electrical Engineering, STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics, Kasteelpark Arenberg 10, B-3001 Leuven, Belgium

Disentanglement is a useful property in representation learning, which increases the interpretability of generative models such as variational autoencoders (VAE), generative adversarial models, and their many variants. Typically in such models, an increase in disentanglement performance is traded off with generation quality. In the context of latent space models, this work presents a representation learning framework that explicitly promotes disentanglement by encouraging orthogonal directions of variations. The proposed objective is the sum of an autoencoder error term along with a principal component analysis reconstruction error in the feature space. This has an interpretation of a restricted kernel machine with the eigenvector matrix valued on the Stiefel manifold. Our analysis shows that such a construction promotes disentanglement by matching the principal directions in the latent space with the directions of orthogonal variation in data space. In an alternating minimization scheme, we use the Cayley ADAM algorithm, a stochastic optimization method on the Stiefel manifold along with the Adam optimizer. Our theoretical discussion and various experiments show that the proposed model is an improvement over many VAE variants in terms of both generation quality and disentangled representation learning.

1 Introduction

Latent space models are popular tools for sampling from high-dimensional distributions. Often, only a small number of latent factors are sufficient to describe data variations. These models exploit the underlying structure of the data and learn explicit representations that are faithful to the data-generating factors. Popular latent space models are variational autoencoders (VAEs; Kingma & Welling, 2014), restricted Boltzmann machines (RBMs; Salakhutdinov & Hinton, 2009), normalizing flows (Rezende & Mohamed, 2015), and their many variants.

In latent variable models, one is often interested in modeling the data in terms of uncorrelated or independent components, yielding a so-called disentangled representation (Bengio, Courville, & Vincent, 2013), which is often studied in the context of VAEs. Generative adversarial networks (GAN) have also been extended to perform disentangled representation learning, for instance, with Info-GANs. It is a GAN that also maximizes the mutual information between a small subset of the discrete latent codes and the true images. In principle, disentanglement corresponds to identifying the underlying factors that generate the data. Components corresponding to the orthogonal directions in latent space may be interpreted as generating distinct factors in the input space (e.g. lighting conditions, style, colors). An illustration of a latent traversal is shown in Figure 1, where one observes that only one specific feature of the image is changing as one moves along a component in the latent space. For instance, in Figure 1, we observe that moving along the first component (vector u_1) generates images where only floor color is varying, while, all other features, such as shape, scale, wall color, and object color, are constant, whereas traversing along the sixth component (vector u_6), for instance, generates images where only the object scale changes as shown in the second row. As we explain later, the components here refer to the principal components given by the principal component analysis (PCA). Therefore, these principal directions encode the directions of maximum variance. Since the floor color is encoded by the largest number of pixels, it gets represented by the first principal component u_1 . Similarly, the other components correspond to the directions with smaller variance. An advantage of such a representation is that the different latent units impart more interpretability to the model. Disentangled models are useful for the generation of plausible pseudo-data with certain desirable properties (e.g., generating new car designs with a predefined color or height).

Now we introduce the mathematical setting to formalize our discussion throughout the paper. We start by introducing a VAE (Kingma & Welling, 2014). Let $p(x)$ be the distribution of the data $x \in \mathbb{R}^d$ and consider latent vectors $z \in \mathbb{R}^\ell$ with the prior distribution $p(z)$, typically a standard normal distribution. Then, one defines an encoder $q(z|x)$ that can be deterministic or probabilistic, for example, given by $\mathcal{N}(z|\phi_\theta(x), \gamma^2\mathbb{I})$, where the

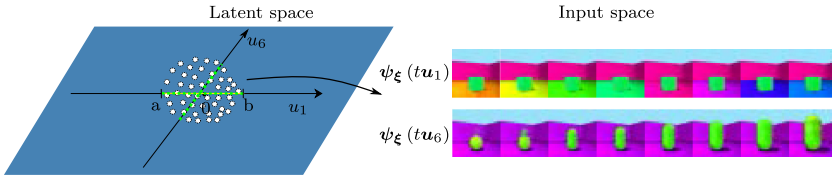


Figure 1: Images by the decoder of the latent space traversal: $\psi_\xi(tu_i)$ for $t \in [a, b]$ with $a < b$ and for some $i \in \{1, \dots, m\}$. Green and black dashed lines represent the walk along u_1 and u_6 , respectively. At every step of the walk, the output of the decoder generates the data in the input space. The images were generated by St-RKM with $\sigma = 10^{-3}$ on 3Dshapes dataset. See Figure 5 for traversal along other components.

mean¹ is given by the neural network ϕ_θ parametrized by θ . A random decoder $p(x|z) = \mathcal{N}(x|\psi_\xi(z), \sigma_0^2\mathbb{I})$ is associated with the decoder neural network ψ_ξ , parameterized by ξ , which maps latent codes to the data points. A VAE is trained by maximizing the lower bound to the idealized log-likelihood as:

$$\mathbb{E}_{z \sim q(z|x)}[\log(p(x|z))] - \beta \text{KL}(q(z|x), p(z)) \leq \log p(x). \tag{1.1}$$

This lower bound is often called as the evidence lower bound (ELBO) when $\beta = 1$. Higgins et al. (2017) show that the larger values of $\beta > 1$ promote more disentanglement but at the expense of generation quality. In this article, we attempt to reconcile the generation quality with disentanglement. To introduce the model, we first make explicit the connection between β -VAEs and standard autoencoders (AEs). Let the data set be $\{x_i\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$. Let $q(z|x) = \mathcal{N}(z|\phi_\theta(x), \gamma^2\mathbb{I})$ be an encoder, where $z \in \mathbb{R}^\ell$. For a fixed $\gamma > 0$, the maximization problem 1.1 is then equivalent to the minimization of the regularized AE,

$$\min_{\theta, \xi} \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}_\epsilon \|x_i - \psi_\xi(\phi_\theta(x_i) + \epsilon)\|_2^2 + \alpha \|\phi_\theta(x_i)\|_2^2 \right\}, \tag{1.2}$$

where $\alpha = \beta\sigma_0^2$, $\epsilon \sim \mathcal{N}(0, \gamma^2\mathbb{I})$ and additive constants depending on γ have been omitted. The first term in equation 1.2 can be interpreted as an AE loss, whereas the second term can be viewed as a regularization. This regularized AE interpretation motivates our method as introduced in section 3.

¹A typical implementation of VAE includes another neural network (after the primary network) for parametrizing the covariance matrix. To simplify this introductory discussion, this matrix is here chosen as a constant diagonal $\gamma^2\mathbb{I}$.

The rest of the article is organized as follows. In section 2 we discuss the closely related work on disentangled representation learning and generation in the context of autoencoders. Further in section 3, we describe the proposed model along with the connection between PCA and disentanglement. In section 3.2, we discuss our contributions. In section 4, we derive the evidence lower bound of the proposed model and show connections with the probabilistic models. In section 5, we describe our experiments and discuss the results.

2 Related Work

Related works can be broadly classified into two categories: Variational autoencoders (VAE) in the context of disentanglement and Restricted Kernel Machines (RKM), a recently proposed modeling framework that integrates kernel methods with deep learning.

2.1 VAE. As discussed in the section 1 (Higgins et al., 2017) suggested that a stronger emphasis on the posterior to match the factorized unit gaussian prior puts further constraints on the implicit capacity of the latent bottleneck. Burgess et al. (2017) further analyzed the effect of the β term in depth. Later, Chen, Li, Grosse, and Duvenaud (2018) showed that the KL term includes the mutual information gap, which encourages disentanglement. Recently, several variants of VAEs promoting disentanglement have been proposed by adding extra terms to the ELBO. For instance, FactorVAE (Kim & Mnih, 2018) augments the ELBO by a new term enforcing factorization of the marginal posterior (or aggregate posterior). Rolínek et al. (2019) analyzed the reason for the alignment of the latent space with the coordinate axes, as the design of VAE itself does not suggest any such mechanism. The authors argue that due to the diagonal approximation in the encoder, together with the inherent stochasticity, forces the local orthogonality of the decoder. Locatello et al. (2020) considered adding an extra term that accounts for the knowledge of some partial label information to improve disentanglement. Later, Ghosh, Sajjadi, Vergari, Black, and Schölkopf (2020) studied the deterministic AEs, where another quadratic regularization on the latent vectors was proposed. In contrast to Rolínek et al. (2019), where the implicit orthogonality of VAE was studied, our proposed model has orthogonality by design due to the introduction of the Stiefel manifold.

2.2 RKM. Restricted kernel machines (RKM; Suykens, 2017) provides a representation of kernel methods with visible and hidden variables similar to the energy function of restricted Boltzmann machines (RBM; LeCun, Huang, & Bottou, 2004; Hinton, 2005), thus linking kernel methods with RBMs. Training and prediction schemes are characterized by the stationary points for the unknowns in the objective. The equations in these stationary points lead to solving a linear-system or matrix decomposition

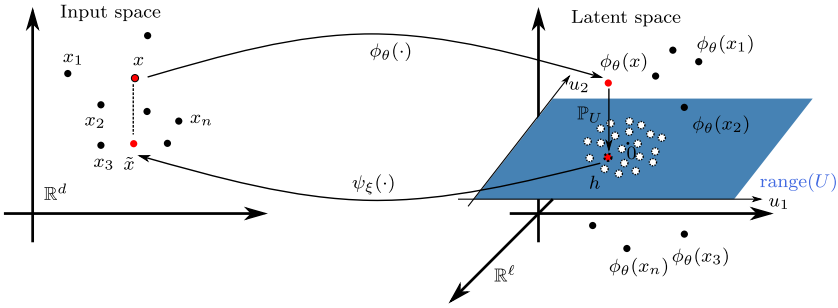


Figure 2: Schematic illustration of St-RKM training problem. The length of the dashed line represents the reconstruction error (see the autoencoder term in equation 3.3) and the length of the vector projecting on hyperplane represents the PCA reconstruction error. After training, the projected points tend to be distributed normally on the hyperplane.

for the training. Suykens (2017) shows various RKM formulations for doing classification, regression, kernel PCA, and singular value decomposition. Later the kernel PCA formulation of RKM was extended to a multiview generative model called generative-RKM (Gen-RKM) which uses convolutional neural networks as explicit feature maps (Pandey, Schreurs, & Suykens, 2020, 2021). For the joint feature selection and subspace learning, the proposed training procedure performs eigendecomposition of the kernel/covariance matrix in every minibatch of the optimization scheme. Intuitively, the model could be seen as learning an autoencoder with kernel PCA in the bottleneck part. As a result, the computational complexity scales cubically with the minibatch size and is proportional to the number of minibatches. Moreover, backpropagation through the eigendecomposition could be numerically unstable due to the possibility of small eigenvalues. All such limitations are addressed by our proposed model.

3 Proposed Mechanism

The main idea of this article consists of learning an autoencoder, along with finding an optimal linear subspace of the latent space such that the variance of the training set in latent space is maximized within this space. (See Figure 2 to follow the discussion below.) Note the distinction with linear autoencoders, which also project the data into the low-dimensional subspace although via nonorthogonal transformations. As a consequence, the latent variables are not guaranteed to be uncorrelated. The encoder $\phi_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^\ell$ typically sends input data to a latent space, while the decoder $\psi_\xi : \mathbb{R}^\ell \rightarrow \mathbb{R}^d$ goes in the reverse direction and constitutes an approximate inverse. Both the encoder and decoder are neural networks parameterized by vectors θ

and ξ . However, it is unclear how to define a parameterization or an architecture of these neural networks so that the learned representation is disentangled. Therefore, in addition to these trained parameters, we also jointly find an m -dimensional linear subspace $\text{range}(U)$ of the latent space \mathbb{R}^ℓ , such that the encoded training points mostly lie within this subspace. This linear subspace is given by the span of the orthonormal columns of the $\ell \times m$ matrix $U = [\mathbf{u}_1, \dots, \mathbf{u}_m]$. The set of such matrices with m orthonormal columns in \mathbb{R}^ℓ with $\ell \geq m$ defines the Stiefel manifold $\text{St}(\ell, m)$. For a reference about optimization on Stiefel manifold, we refer to Absil, Mahony, and Sepulchre (2008). Input data are then encoded into a subspace of the latent space by

$$\mathbf{x} \mapsto \mathbb{P}_U \boldsymbol{\phi}_\theta(\mathbf{x}) = \mathbf{u}_1^\top \boldsymbol{\phi}_\theta(\mathbf{x}) \times \begin{bmatrix} | \\ \mathbf{u}_1 \\ | \end{bmatrix} + \dots + \mathbf{u}_m^\top \boldsymbol{\phi}_\theta(\mathbf{x}) \times \begin{bmatrix} | \\ \mathbf{u}_m \\ | \end{bmatrix},$$

where the orthogonal projector onto $\text{range}(U)$ is simply $\mathbb{P}_U = UU^\top$.

Orthogonal latent directions. Naturally, given an $m \times m$ orthogonal matrix O and a matrix $U \in \text{St}(\ell, m)$, we have

$$\text{range}(U) = \text{range}(UO).$$

To select a specific matrix $U_\star = [\mathbf{u}_{\star,1}, \dots, \mathbf{u}_{\star,m}] \in \text{St}(\ell, m)$, we choose $\mathbf{u}_{\star,1}, \dots, \mathbf{u}_{\star,m}$ to be the eigenvectors of the matrix $C_\theta = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\phi}_\theta(\mathbf{x}_i) \boldsymbol{\phi}_\theta^\top(\mathbf{x}_i)$, associated with the m largest eigenvalues sorted in descending order. For simplicity, we assume that the m largest eigenvalues of C_θ are distinct, whereas the general case involves minor technicalities. Here the feature map is assumed to be centered, $\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})}[\boldsymbol{\phi}_\theta(\mathbf{x})] = \mathbf{0}$, so that C_θ is interpreted as a covariance matrix. Next, we state a result that we will use extensively later.

Proposition 1. *Let M be an $\ell \times \ell$ symmetric matrix. Let v_1, \dots, v_m be its m smallest eigenvalues, possibly including multiplicities, with associated orthonormal eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_m$. Let V be a matrix whose columns are these eigenvectors. Then the optimization problem $\min_{U \in \text{St}(\ell, m)} \text{Tr}(U^\top M U)$ has a minimizer at $U_\star = V$ and we have $U_\star^\top M U_\star = \text{diag}(\mathbf{v})$, with $\mathbf{v} = (v_1, \dots, v_m)^\top$.*

A few remarks follow. First, if U_\star is a minimizer of the optimization problem in proposition 1 then $U'_\star = U_\star O$ with O orthogonal is also a minimizer, but $U'^\top_\star M U'_\star$ is not necessarily diagonal. Second, notice that if the eigenvalues of M in proposition 1 have a multiplicity larger than 1, there can exist several sets of eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_m$, associated with the m smallest eigenvalues, spanning distinct linear subspaces. Nevertheless, in practice, the eigenvalues of the matrices considered in this article are numerically distinct.

We now use proposition 1. For a given positive integer $m \leq \ell$, the subspace spanned by the eigenvectors of C_θ with the m largest eigenvalues is obtained by solving

$$\min_{U \in \text{St}(\ell, m)} \text{Tr}(C_\theta - \mathbb{P}_U C_\theta \mathbb{P}_U) = \frac{1}{n} \sum_{i=1}^n \|\mathbb{P}_{U^\perp} \phi_\theta(x_i)\|_2^2,$$

where $\mathbb{P}_{U^\perp} = \mathbb{I} - \mathbb{P}_U$, as it is explained, for instance, in section 4.1 of Avron, Nguyen, and Woodruff (2014). The above objective corresponds to the reconstruction error of kernel PCA, for the kernel $k_\theta(x, y) = \phi_\theta^\top(x)\phi_\theta(y)$. As described earlier, we choose a specific $U_\star \in \text{St}(\ell, m)$ by requiring that the following matrix is diagonal,

$$U_\star^\top C_\theta U_\star = \text{diag}(\lambda), \tag{3.1}$$

where λ is a vector containing the m largest eigenvalues sorted in decreasing order. If these eigenvalues are distinct, then the U_\star is essentially unique, up to sign flip of each of its columns. Notice that $\text{Tr}(U_\star^\top C_\theta U_\star) = \text{Tr}(U_\star U_\star^\top C_\theta U_\star U_\star^\top)$.

Orthogonal directions of variation in input space. We want the lines defined by the orthonormal vectors $\{u_{\star,1}, \dots, u_{\star,m}\}$ to provide directions associated with different generative factors of our model. In other words, we conjecture that a possible formalization of disentanglement is that the principal directions in latent space match orthogonal directions of variation in the data space (see Figure 2). That is, we would like that

$$U_\star^\top \sum_{a=1}^d (\nabla \psi_a(y_i) \nabla \psi_a(y_i)^\top) U_\star \text{ is diagonal}, \tag{3.2}$$

for all the points in latent space $y_i = \mathbb{P}_U \phi_\theta(x_i)$ for $i = 1, \dots, n$. In equation 3.2, $\psi_a(y)$ refers to the a th component of the image $\psi(y) \in \mathbb{R}^d$. To sketch this idea, we study the local motions in the latent space.

Let $\Delta_k = \nabla \psi(y)^\top u_{\star,k} \in \mathbb{R}^d$ be the directional derivative of ψ at point y in the direction $u_{\star,k}$ with $1 \leq k \leq m$. Then, as one moves in the latent space from a point y in the direction of $u_{\star,k}$, the generated data change by

$$\psi(y + t u_{\star,k}) - \psi(y) = t \Delta_k + \mathcal{O}(t^2),$$

with $\Delta_k \in \mathbb{R}^d$ and $t \in \mathbb{R}$. Consider now a different direction, $k' \neq k$. As the latent point moves along $u_{\star,k}$ or along $u_{\star,k'}$, we expect the decoder output to vary in a significantly different manner, $\Delta_k^\top \Delta_{k'} = 0$. We presume this interpretation to model the change in floor color and object scale in Figure 1 for instance. More explicitly, we can expect u_k and $u_{k'}$ to model, respectively, the change of colors of the floor and of the main object while leaving the color of the other objects unchanged. Since the floor and the main object

do not overlap, that is, they are different regions in pixel space, we would have $\Delta_k^\top \Delta_k = 0$. Admittedly, the change in object shape in Figure 1 is less obviously interpreted. Now, denote by Δ the matrix obtained by stacking the vector Δ_k as columns for $1 \leq k \leq m$. Explicitly, we have $\Delta = \nabla \psi_a(\mathbf{y})^\top U_\star$. Hence, for all \mathbf{y} in the latent space, we expect the Gram matrix $\Delta^\top \Delta$ to be diagonal (see equation 3.2). We now discuss how this idea might be realized by minimizing specific objective functions.

3.1 Objective Function. In this article, we propose to train an objective function which is composed of an AE loss and a PCA loss. Hence, the proposed model is given by

$$\min_{\substack{U \in \text{St}(\ell, m) \\ \theta, \xi}} \lambda \underbrace{\frac{1}{n} \sum_{i=1}^n L_{\xi, \mathbb{P}_U}(x_i, \phi_\theta(x_i))}_{\text{Autoencoder objective}} + \underbrace{\text{Tr}(C_\theta - \mathbb{P}_U C_\theta \mathbb{P}_U)}_{\text{PCA objective}}, \quad (3.3)$$

where $\lambda > 0$ is a trade-off parameter and $C_\theta = \frac{1}{n} \sum_{i=1}^n \phi_\theta(x_i) \phi_\theta^\top(x_i)$. Naturally, the above objective is invariant if U is replaced by UO with O an orthogonal matrix. Given a local minimizer, we select $U_\star \in \text{St}(\ell, m)$ such that $U_\star^\top C_\theta U_\star$ is diagonal as in equation 3.1, to identify the principal directions in the latent space. This last step is conveniently done with a singular value decomposition (see step 10 of algorithm 1). In the proposed model, reconstruction of an out-of-sample point x is given by $\psi_\xi(\mathbb{P}_U \phi_\theta(x))$. We call the procedure to

find a triplet (U_\star, θ, ξ) solving (5) s.t. $U_\star^\top C_\theta U_\star$ is diagonal, St-RKM

the training of a Stiefel-restricted kernel machines, equation 3.3, in view of our discussion in section 2. The basic idea is to design different AE losses with a regularization term that penalizes the feature map in the orthogonal subspace U^\perp . The choice of the AE losses is motivated by the expression of the regularized AE in equation 1.2 and by the following lemma, which extends the result of Rolínek et al. (2019). Here we adapt it in the context of optimization on the Stiefel manifold (see appendix for the proof).

Lemma 1. *Let $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_m)$ a random vector and $U \in \text{St}(\ell, m)$. Let $\psi_a(\cdot) \in C^2(\mathbb{R}^\ell)$ with $a \in [d]$. If the function $[\psi(\cdot) - x]_a^2$ has L_a -Lipschitz continuous Hessian for all $a \in [d]$, we have*

$$\begin{aligned} \mathbb{E}_\epsilon \|x - \psi(\mathbf{y} + \sigma U \epsilon)\|_2^2 &= \|x - \psi(\mathbf{y})\|_2^2 + \sigma^2 \text{Tr}(U^\top \nabla \psi(\mathbf{y}) \nabla \psi(\mathbf{y})^\top U) \\ &\quad - \sigma^2 \sum_{a=1}^d [x - \psi(\mathbf{y})]_a \text{Tr}(U^\top \text{Hess}_y[\psi_a] U) + \sum_{a=1}^d R_a(\sigma), \end{aligned} \quad (3.4)$$

with $|R_a(\sigma)| \leq \frac{1}{6} \sigma^3 L_a \frac{\sqrt{2(m+1)} \Gamma((m+1)/2)}{\Gamma(m/2)}$ where Γ is Euler's gamma function.

In lemma 1, the first term on the right-hand side in equation 3.4 plays the role of the classical AE loss. The second term is proportional to the trace of equation 3.2. This is related to our discussion above where we argue that jointly diagonalizing both $U^\top \nabla \psi(\mathbf{y}) \nabla \psi(\mathbf{y})^\top U$ and $U^\top C_\theta U$ helps to enforce disentanglement. However, determining the behavior of the third term in equation 3.4 is difficult. This is because, for a typical neural network architecture, it is unclear in practice if the function $\|x - \psi(\cdot)\|_a^2$ has L_a -Lipschitz continuous Hessian for all $a \in [d]$. Hence we propose another AE loss (split loss) in order to cancel the third term in equation 3.4. Nevertheless, the assumption in lemma 1 is used to provide a meaningful bound on the remainder in equation 3.4. In the light of these remarks, we propose two stochastic AE losses.

3.1.1 AE Losses. In analogy with the VAE objective equation 1.2, the first AE encoder loss function can be chosen as

$$L_{\xi, \mathbb{P}_U}^{(\sigma)}(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbb{I}_m)} \|\mathbf{x} - \psi_\xi(\mathbb{P}_U \mathbf{z} + \sigma U \epsilon)\|_2^2, \text{ with } \sigma > 0.$$

As motivated by lemma 1, the noise term $\sigma U \epsilon$ above promotes a smoother decoder network. To further promote disentanglement, we propose a split AE loss

$$L_{\xi, \mathbb{P}_U}^{(\sigma), sl}(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \psi_\xi(\mathbb{P}_U \mathbf{z})\|_2^2 + \mathbb{E}_\epsilon \|\psi_\xi(\mathbb{P}_U \mathbf{z}) - \psi_\xi(\mathbb{P}_U \mathbf{z} + \sigma U \epsilon)\|_2^2, \quad (3.5)$$

with $\epsilon \sim \mathcal{N}(0, \mathbb{I}_m)$. The first term in equation 3.5 is the classical AE loss while the second term promotes orthogonal directions of variations. Thus, by relating lemma 1 to equation 3.5 we see that

$$L_{\xi, \mathbb{P}_U}^{(\sigma), sl}(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \psi_\xi(\mathbb{P}_U \mathbf{z})\|_2^2 + \sigma^2 \text{Tr}(U^\top \nabla \psi(\mathbf{y}) \nabla \psi(\mathbf{y})^\top U) + \sum_{a=1}^d R_a(\sigma).$$

In short, the optimization over U in equation 3.3 with the splitted loss aims to promote a U_\star such that

$$U_\star^\top C_\theta U_\star \text{ and } U_\star^\top \left(\sum_{i=1}^n \nabla \psi(\mathbf{y}_i) \nabla \psi(\mathbf{y}_i)^\top \right) U_\star \text{ are jointly diagonal.}$$

Figure 3 gives a visualization of the diagonal form of

$$\frac{1}{|C|} \sum_{i \in C} U_\star^\top \nabla \psi(\mathbf{y}_i) \nabla \psi(\mathbf{y}_i)^\top U_\star, \text{ with } \mathbf{y}_i = \mathbb{P}_U \phi_\theta(\mathbf{x}_i) \quad (3.6)$$

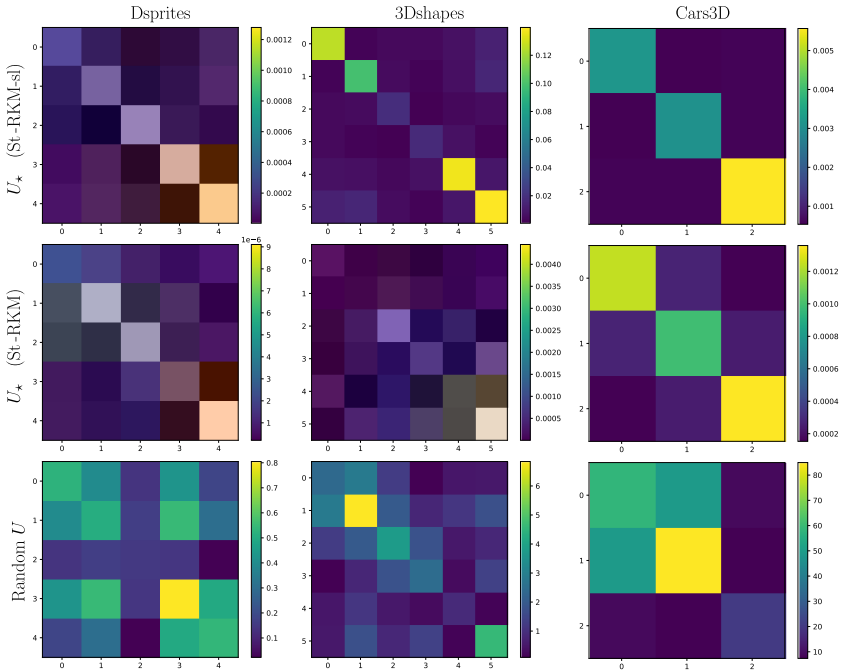


Figure 3: Visualizing the matrix, equation 3.6 for St-RKM models after training on three data sets. The first two rows show, equation 3.6, where $U = U_* \in \text{St}(\ell, m)$ is the output of algorithm 1. These matrices are effectively close to being diagonal and especially for St-RKM-sl, as expected. In contrast, the third row shows the same matrix, equation 3.6, with $U \in \text{St}(\ell, m)$ sampled uniformly at random (see Table 6 for the corresponding normalized diagonalization errors).

obtained after training; where \mathcal{C} contains the indices of a subset of 50 images sampled uniformly at random. (For numerical values, Table 6 in the appendix shows the normalized diagonalization errors.)

Note that we do not simply propose another encoder-decoder architecture, given by $U^\top \phi_\theta(\cdot)$ and $\psi_\xi(U \cdot)$. Instead, our objective assumes that the neural network defining the encoder provides a better embedding if we impose that it maps training points on a linear subspace of dimension $m < \ell$ in the ℓ -dimensional latent space. In other words, the optimization of the parameters in the last layer of the encoder does not play a redundant role, since the second term in equation 3.3 clearly also depends on $\mathbb{P}_{U^\perp} \phi_\theta(\cdot)$. The full training involves an alternating minimization procedure, which is described in algorithm 1.

3.2 Contributions. Here is a summary of our contributions. We propose three main changes with respect to the related works. First, to promote disentangled representation learning, we propose orthogonal projection in the latent space via a rectangular matrix that is valued on the Stiefel manifold. Then for the training, we use the Cayley ADAM algorithm of Li, Li, and Todorovic (2020) for stochastic optimization on the Stiefel manifold and call our proposed model St-RKM. Second, we propose several objective functions to learn the feature map and the pre-image map networks in the form of an encoder and a decoder, respectively. The best configuration for promoting a disentangled representation is

$$\min_{\substack{U \in \text{St}(\ell, m) \\ \theta, \xi}} \frac{\lambda}{n} \sum_{i=1}^n (\text{splitted}) \text{ AE loss}(x_i, \mathbb{P}_U, \theta, \xi) + \text{PCA objective}(C_\theta, \mathbb{P}_U),$$

where the covariance matrix reads $C_\theta = \frac{1}{n} \sum_{i=1}^n \phi_\theta(x_i)\phi_\theta^\top(x_i)$ and $\mathbb{P}_U = UU^\top$ with U an $\ell \times m$ matrix with orthonormal columns. Here $\lambda > 0$ is a trade-off parameter. The final parameters (U_\star, θ, ξ) give a local minimizer of this objective with U_\star chosen such that $U_\star^\top C_\theta U_\star$ is diagonal. Third, we validate through experiments the following statement: The combination of a split AE loss with a PCA objective by using an explicit optimization on the Stiefel manifold promotes disentanglement. In this article, disentanglement is interpreted as jointly diagonalizing the matrix representing variations in the input space with respect to latent motions $\sum_i U_\star^\top \nabla \psi_\xi(y_i) \nabla \psi_\xi(y_i)^\top U_\star$ where $y_i = \mathbb{P}_U \phi_\theta(x_i)$ and the covariance matrix of the data set in the latent space $U_\star^\top C_\theta U_\star$.

4 Connections with the Evidence Lower Bound

We now discuss the interpretation of the proposed model in the probabilistic setting and the independence of latent factors. In order to formulate an ELBO, consider the following random encoders,

$$q(z|x) = \mathcal{N}(z|\phi_\theta(x), \gamma^2 \mathbb{I}_\ell) \text{ and } q_U(z|x) = \mathcal{N}(z|\mathbb{P}_U \phi_\theta(x), \sigma^2 \mathbb{P}_U + \delta^2 \mathbb{P}_{U^\perp}),$$

where ϕ_θ has zero mean on the data distribution. Here, σ^2 plays the role of a trade-off parameter, while the regularization parameter δ is introduced for technical reasons and is put to a numerically small absolute value (see the appendix for details). Let the decoder be $p(x|z) = \mathcal{N}(x|\psi_\xi(z), \sigma_0^2 \mathbb{I})$ and the latent space distribution is parameterized by $p(z) = \mathcal{N}(0, \Sigma)$ where $\Sigma \in \mathbb{R}^{\ell \times \ell}$ is a covariance matrix. We treat Σ as a parameter of the optimization problem that is determined at the last stage of the training. Then the minimization problem 3.3 with stochastic AE loss is equivalent to the

maximization of

$$\frac{1}{n} \sum_{i=1}^n \left\{ \underbrace{\mathbb{E}_{q_U(z|x_i)}[\log(p(x_i|z))]}_{(I)} - \underbrace{\text{KL}(q_U(z|x_i), q(z|x_i))}_{(II)} - \underbrace{\text{KL}(q_U(z|x_i), p(z))}_{(III)} \right\}, \quad (4.1)$$

which is a lower bound to the ELBO, since the KL divergence in term II in equation 4.1 is positive. For details of the derivation, see the appendix. The hyperparameters γ , σ , σ_0 take a fixed value. Up to additive constants, the terms I and II of equation 4.1 match the objective, equation 3.3. The third term (III) in equation 4.1 is optimized after the training of the first two terms. It can be written as

$$\frac{1}{n} \sum_{i=1}^n \text{KL}(q_U(z|x_i), p(z)) = \frac{1}{2} \text{Tr}[\Sigma_0 \Sigma^{-1}] + \frac{1}{2} \log(\det \Sigma) + \text{constants},$$

with $\Sigma_0 = \mathbb{P}_U C_\theta \mathbb{P}_U + \sigma^2 \mathbb{P}_U + \delta^2 \mathbb{P}_{U^\perp}$. In that case, the optimal covariance matrix is diagonalized $\Sigma = U(\text{diag}(\lambda) + \sigma^2 \mathbb{I}_m)U^\top + \delta^2 \mathbb{P}_{U^\perp}$, with λ denoting the principal values of the PCA.

Now we briefly discuss the factorization of the encoder. Let $h(x) = U^\top \phi_\theta(x)$ and let the effective latent variable be $z^{(U)} = U^\top z \in \mathbb{R}^m$. Then the probability density function of $q_U(z|x)$ is

$$f_{q_U(z|x)}(z) = \frac{e^{-\frac{\|U^\top z\|_2^2}{2\delta^2}}}{(\sqrt{2\pi}\delta^2)^{\ell-m}} \prod_{j=1}^m \frac{e^{-\frac{(z_j^{(U)} - h_j(x))^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma^2},$$

where the first factor is approximated by a Dirac delta if $\delta \rightarrow 0$. Hence, the factorized form of q_U shows the independence of the latent variables $z^{(U)}$. This factorization is used as a regularization term in the objective by Kim and Mnih (2018) to promote disentanglement. In particular, term II in equation 4.1 is analogous to a “total correlation” loss (Chen et al., 2018).

5 Experiments

In this section, we investigate if St-RKM² can simultaneously achieve accurate reconstructions on training data, good random generations, and good disentanglement performance. We use the standard data sets: MNIST (LeCun & Cortes, 2010), Fashion-MNIST (fMNIST; Xiao, Rasul, & Vollgraf, 2017), and SVHN (Netzer et al., 2011). To evaluate disentanglement, we use data sets with known ground-truth generating factors such as dSprites

²The source code is available at http://bit.ly/StRKM_code.

Algorithm 1: Manifold Optimization of St-RKM.

Input: $\{\mathbf{x}_i\}_{i=1}^n, \phi_\theta, \psi_\zeta, \mathcal{J} := \text{equation 3.3}$
Output: Learned θ, ζ, U

```

1: procedure TRAIN
2:   while not converged do
3:      $\{\mathbf{x}\} \leftarrow \{\text{Get minibatch}\}$ 
4:     Get embeddings  $\phi_\theta(\mathbf{x}) \leftarrow \mathbf{x}$ 
5:     Compute centered  $C_\theta$  ▷ Covariance matrix
6:     Update  $\{\theta_e, \psi_g\} \leftarrow \text{Adam}(\mathcal{J})$  ▷ Optimization step
7:     Update  $\{U\} \leftarrow \text{Cayley\_Adam}(\mathcal{J})$  ▷ Optimization step
8:   end while
9:   Do steps 4-5 over whole data set
10:   $U \leftarrow \text{SVD}(C_\theta)$  ▷ Equation (3.1)
11: end procedure

```

(Matthey, Higgins, Hassabis, & Lerchner, 2017), 3DShapes (Burgess & Kim, 2018), and 3D cars (Reed, Zhang, Zhang, & Lee, 2015). Further, all figures and tables report average errors with 1 standard deviation over 10 experiments.

5.1 Algorithm. We use an alternating-minimization scheme as shown in algorithm 1. First, the Adam optimizer with a learning rate 2×10^{-4} is used to update the encoder-decoder parameters; then, the Cayley Adam optimizer (Li et al., 2020) with a learning rate 10^{-4} is used to update U . Finally, at the end of the training, we recompute U from the singular value decomposition (SVD) of the covariance matrix as a final correction-step of the kernel PCA term in our objective (step 10 of algorithm 1). Since the $\ell \times \ell$ covariance matrix is typically small, this decomposition is fast (see Table 3). In practice, our training procedure only marginally increases the computation cost, which can be seen from training times in Table 1.

5.2 Experimental Setup. We consider four baselines for comparison: VAE, β -VAE, FactorVAE, and Info-GAN. An ablation study with the

Table 1: Training Time in Minutes (for 1000 Epochs, Mean with 1 Standard Deviation over 10 Runs) and the Number of Parameters (Nb) of the Generative Models on the MNIST Data Set.

Model	St-RKM	(β)-VAE	FactorVAE	Info-GAN
Nb parameters	4164519	4165589	8182591	4713478
Training time	21.93 (1.3)	19.83 (0.8)	33.31 (2.7)	45.96 (1.6)

Gen-RKM is shown in section A.4 in the appendix. Extensive experimentation was not computationally feasible since the evaluation and decomposition of kernel matrices scales $\mathcal{O}(n^2)$ and $\mathcal{O}(n^3)$ with the data set size (see the discussion in section 2).

5.3 Inductive Biases. To be consistent in evaluation, we keep the same encoder (discriminator) and decoder (generator) architecture and the same latent dimension across the models. We use convolutional neural networks due to the choice of image data sets for evaluating generation and disentanglement. In the case of Info-GAN, batch normalization is added for training stability (see section A.3 in the appendix for details). For the determination of the hyperparameters of other models, we start from values in the range of the parameters suggested in the authors’ reference implementation. After trying various values, we noticed that $\beta = 3$ and $\gamma = 12$ seem to work well across the data sets that we considered for β -VAE and FactorVAE, respectively. Furthermore, in all the experiments on St-RKM, we keep the reconstruction weight $\lambda = 1$. All models are trained on the entire data set. Note that for the same encoder-decoder network, the St-RKM model has the least number of parameters compared to any VAE variants and Info-GAN (see Table 1).

To evaluate the quality of generated samples, we report the Fréchet inception distance (FID; Heusel et al., 2017) and the sliced Wasserstein distance (SWD; Karras, Aila, Laine, & Lehtinen, 2017) scores with mean and standard deviation in Figure 4. Note that FID scores are not necessarily appropriate for dSprites since this data set is significantly different from ImageNet on which the Inception network was originally trained. (Randomly generated samples are shown in Figure 8 in the appendix). To generate samples from the deterministic St-RKM ($\sigma = 0$), we sample from a fitted normal distribution on the latent embedding of the data set; for a similar procedure, see Ghosh et al., 2020). Figure 4 shows that the St-RKM variants perform better (lower mean scores) on most data sets, and within them, the stochastic variants with $\sigma = 10^{-3}$ perform best. This can be attributed to a better generalization of the decoder network due to the addition of noise term on latent variables (see lemma 1). The training times for St-RKM variants are shorter compared to FactorVAE and Info-GAN due to a significantly small number of parameters.

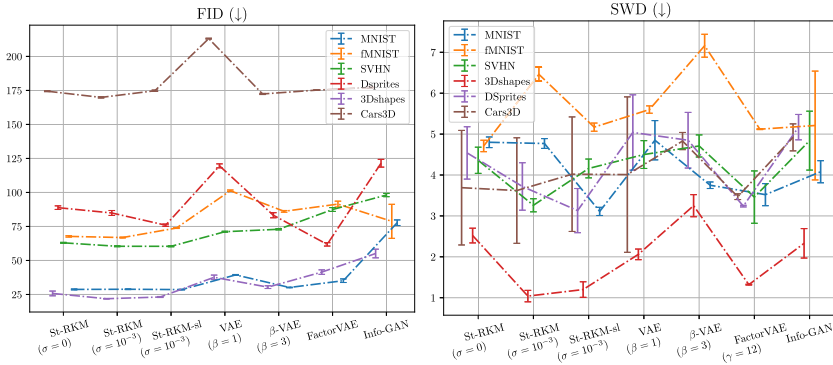


Figure 4: Fréchet inception distance (FID; Heusel, Ramsauer, Unterthiner, Nessler, & Hochreiter, 2017) and sliced Wasserstein distance (SWD scores (mean and 1 standard deviation) for 8000 randomly generated samples (smaller is better).

To evaluate the disentanglement performance, various metrics have been proposed. A comprehensive review by Locatello et al. (2019) shows that the various disentanglement metrics are correlated, albeit with a different degree of correlation across data sets. In this article, we use three metrics to evaluate disentanglement: Eastwood’s framework (Eastwood & Williams, 2018), mutual information gap (MIG; Chen et al., 2018), and separated attribute predictability (SAP; Kumar et al., 2018) scores. Eastwood’s framework (Eastwood & Williams, 2018) further proposes three metrics: *disentanglement*: the degree to which a representation factorizes the underlying factors of variation, with each variable capturing at most one generative factor; *completeness*: the degree to which each underlying factor is captured by a single code variable; and *informativeness*: the amount of information that a representation captures about the underlying factors of variation. Furthermore, we use a slightly modified version of MIG score as proposed by Locatello et al. (2019). Figure 6 shows that St-RKM variants have better disentanglement and completeness scores (higher mean scores). However, the informativeness scores are higher for St-RKM when using a lasso-regressor in contrast to mixed scores with a random forest regressor. Figure 7 further complements these observations by showing MIG and SAP scores. Here, the St-RKM-sl model has the highest mean scores for every data set. Qualitative assessment can be done from Figure 5, which shows the generated images by traversing along the principal components in the latent space. In the 3DShapes data set, the St-RKM model captures floor hue, wall hue, and orientation perfectly but has a slight entanglement in capturing other factors. This is worse in β -VAE, which has entanglement in all dimensions except the floor hue, along with noise in some generated images. Similar trends can be observed in the dSprites and 3D cars data sets.

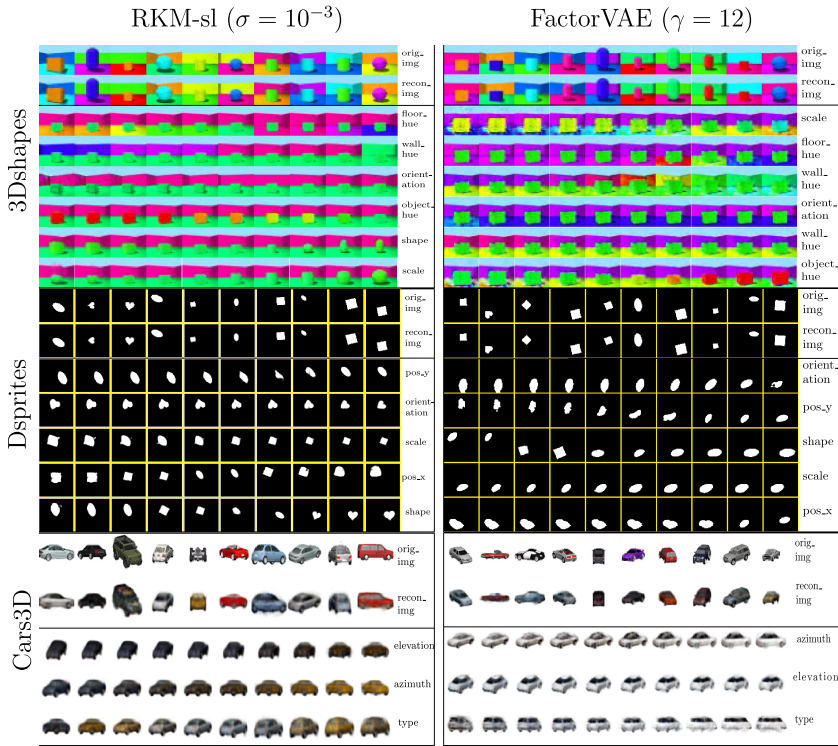


Figure 5: Traversals along the principal components. The first two rows show the ground-truth and reconstructed images. Each subsequent row shows the generated images by traversing along a principal component in the latent space. The last column in each subimage indicates the dominant factor of variation.

6 Conclusion

This article proposes the St-RKM model for disentangled representation learning and generation based on manifold optimization. For the training, we use the Cayley Adam algorithm of Li et al. (2020) for stochastic optimization on the Stiefel manifold. Computationally, St-RKM increases the training time by only a reasonably small amount compared to β -VAE, for instance. Furthermore, we propose several autoencoder objectives and discuss that the combination of a stochastic AE loss with an explicit optimization on the Stiefel manifold promotes disentanglement. In addition, we establish connections with probabilistic models, formulate an evidence lower bound, and discuss the independence of latent factors. Where the considered baselines have a trade-off between generation quality and disentanglement, we improve on both of these aspects as illustrated through

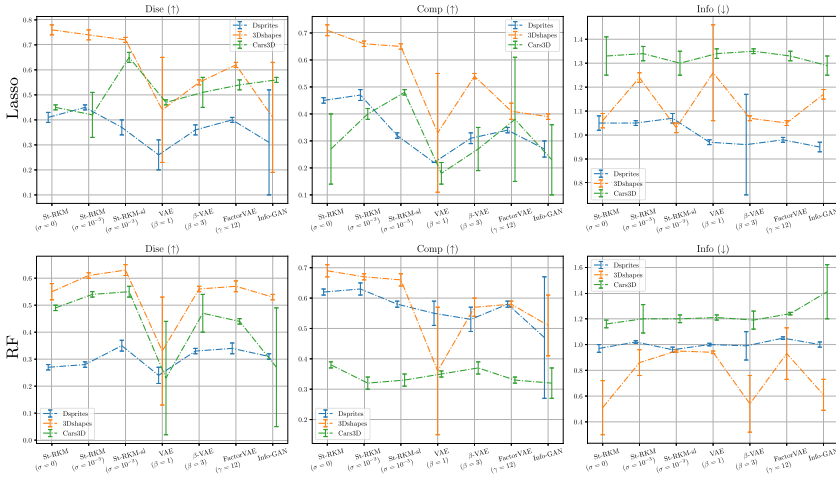


Figure 6: Eastwood framework’s (Eastwood & Williams, 2018) disentanglement metric with Lasso and random forest (RF) regressor. The plot shows mean and 1 standard deviation of scores over 10 iterations. For disentanglement and completeness, a higher score is better; for informativeness, lower is better. “Info.” indicates (average) root-mean-square error in predicting z .

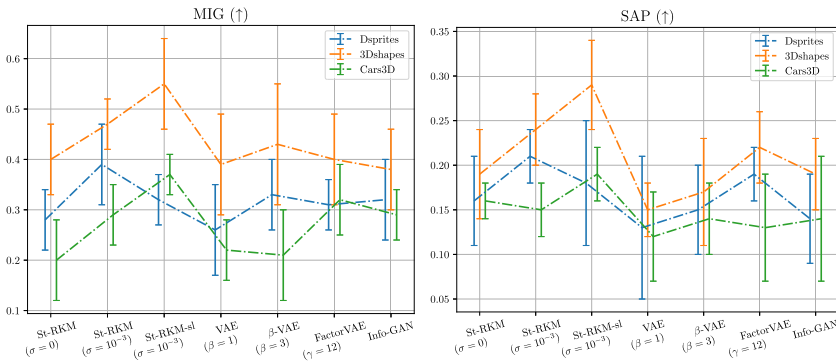


Figure 7: MIG (Chen et al., 2018; Locatello et al., 2019) and SAP (Kumar, Sattigeri, & Balakrishnan, 2018) scores to evaluate disentanglement performance showing the mean (standard deviation) over 10 random seeds.

various experiments. The proposed model has some limitations. A first limitation is hyperparameter selection: the number of components in the KPCA, neural network architecture, and the final size of the feature map. When additional knowledge on the data is available, we suggest that the user selects the number of components close to the number of underlying generating

factors. The final size of the feature map should be large enough so that KPCA extracts meaningful components. Second, we interpret the disentanglement as the two orthogonal changes in the latent space corresponding to two orthogonal changes in input space. Although not perfect, we believe it is a reasonable mathematical approximation of the loosely defined notion of disentanglement. Moreover, experimental results confirm this assumption. Among the possible regularizers on the hidden features, the model associated with the squared Euclidean norm was analyzed in detail, while a deeper study of other regularizers is a prospect for further research, in particular for the case of spherical units.

Appendix

A.1 Proof of Lemma 1. We first quote a result that is used in the context of optimization (Nesterov, 2014, lemma 1.2.4). Let f be a function with L_a -Lipschitz continuous Hessian. Then,

$$\begin{aligned} & \left| \underbrace{f(\mathbf{y}_1) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{y}_1 - \mathbf{y}) - \frac{1}{2} (\mathbf{y}_1 - \mathbf{y})^\top \text{Hess}_y[f](\mathbf{y}_1 - \mathbf{y})}_{r(\mathbf{y}_1 - \mathbf{y})} \right| \\ & \leq \frac{L_a}{6} \|\mathbf{y}_1 - \mathbf{y}\|_2^3. \end{aligned} \quad (\text{A.1})$$

Then we calculate the power series expansion of $f(\mathbf{y}) = [\mathbf{x} - \boldsymbol{\psi}(\mathbf{y})]_a^2$ and take the expectation with respect to $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbb{I})$. First, we have $\nabla f(\mathbf{y}) = -2[\mathbf{x} - \boldsymbol{\psi}(\mathbf{y})]_a \nabla \boldsymbol{\psi}_a(\mathbf{y})$ and

$$\text{Hess}_y[f] = 2\nabla \boldsymbol{\psi}_a(\mathbf{y}) \nabla \boldsymbol{\psi}_a(\mathbf{y})^\top - 2[\mathbf{x} - \boldsymbol{\psi}(\mathbf{y})]_a \text{Hess}_y[\boldsymbol{\psi}_a].$$

Then we use equation A.1 with $\mathbf{y}_1 - \mathbf{y} = \sigma U \boldsymbol{\epsilon}$. By taking the expectation over $\boldsymbol{\epsilon}$, notice that the order 1 term in σ vanishes since $\mathbb{E}_\boldsymbol{\epsilon}[\boldsymbol{\epsilon}] = 0$. We find

$$\begin{aligned} \mathbb{E}_\boldsymbol{\epsilon}[\mathbf{x} - \boldsymbol{\psi}(\mathbf{y} + \sigma U \boldsymbol{\epsilon})]_a^2 &= [\mathbf{x} - \boldsymbol{\psi}(\mathbf{y})]_a^2 + \sigma^2 \text{Tr}(U^\top \nabla \boldsymbol{\psi}_a(\mathbf{y}) \nabla \boldsymbol{\psi}_a(\mathbf{y})^\top U) \\ &\quad - \sigma^2 [\mathbf{x} - \boldsymbol{\psi}(\mathbf{y})]_a \text{Tr}(U^\top \text{Hess}_y[\boldsymbol{\psi}_a] U) + \mathbb{E}_\boldsymbol{\epsilon} r(\sigma U \boldsymbol{\epsilon}), \end{aligned}$$

where we used that $\mathbb{E}_\boldsymbol{\epsilon}[\boldsymbol{\epsilon}^\top M \boldsymbol{\epsilon}] = \text{Tr}[M]$ for any symmetric matrix M since $\mathbb{E}_\boldsymbol{\epsilon}[\boldsymbol{\epsilon}_i \boldsymbol{\epsilon}_j] = \delta_{ij}$. Next, denote $R_a(\sigma) = \mathbb{E}_\boldsymbol{\epsilon} r(\sigma U \boldsymbol{\epsilon})$; we can use the Jensen inequality and subsequently equation A.1:

$$|R_a(\sigma)| = |\mathbb{E}_\boldsymbol{\epsilon} r(\sigma U \boldsymbol{\epsilon})| \leq \mathbb{E}_\boldsymbol{\epsilon} |r(\sigma U \boldsymbol{\epsilon})| \leq \frac{L_a}{6} \mathbb{E}_\boldsymbol{\epsilon} \|\sigma U \boldsymbol{\epsilon}\|_2^3.$$

Next, we notice that $\|\sigma U \epsilon\|_2 = \sigma (\epsilon^\top U^\top U \epsilon)^{1/2} = \sigma \|\epsilon\|_2$. It is useful to notice that $\|\epsilon\|_2$ is distributed according to a chi distribution. By using this remark, we find

$$|R_a(\sigma)| \leq \sigma^3 \frac{L_a}{6} \mathbb{E}_\epsilon \|\epsilon\|_2^3 = \sigma^3 \frac{L_a}{6} \frac{\sqrt{2}(m+1)\Gamma((m+1)/2)}{\Gamma(m/2)},$$

where the last equality uses the expression for the third moment of the chi distribution and where the gamma function Γ is the extension of the factorial to the complex numbers.

A.2 Details on Evidence Lower Bound for St-RKM model. Now we discuss the details of ELBO given in section 4. The first term in equation 4.1 is

$$\begin{aligned} \mathbb{E}_{q_U(z|x_i)}[\log(p(x_i|z))] &= -\frac{1}{2\sigma_0^2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbb{I})} \|x_i - \psi_\xi(\mathbb{P}_U \phi_\theta(x_i) + \sigma \mathbb{P}_U \epsilon + \delta \mathbb{P}_{U^\perp} \epsilon)\|_2^2 \\ &\quad - \frac{d}{2} \log(2\pi \sigma_0^2), \end{aligned}$$

where we used the following reparameterization following Kingma and Welling (2014): $\mathbb{E}_{q_U(z|x_i)}[f(z)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbb{I})} [f(\mathbb{P}_U \phi_\theta(x) + (\sigma \mathbb{P}_U + \delta \mathbb{P}_{U^\perp}) \epsilon)]$, with $p(x|z) = \mathcal{N}(x | \psi_\xi(z), \sigma_0^2 \mathbb{I})$, and $q_U(z|x) = \mathcal{N}(z | \mathbb{P}_U \phi_\theta(x), \sigma^2 \mathbb{P}_U + \delta^2 \mathbb{P}_{U^\perp})$. Clearly, the above expectation can be written as

$$\mathbb{E}_\epsilon \mathbb{E}_{\epsilon_\perp} \|x_i - \psi_\xi(\mathbb{P}_U \phi_\theta(x_i) + \sigma U \epsilon + \delta U_\perp \epsilon_\perp)\|_2^2,$$

with $\epsilon \sim \mathcal{N}(0, \mathbb{I}_m)$ and $\epsilon_\perp \sim \mathcal{N}(0, \mathbb{I}_{\ell-m})$. Hence, we fix $\sigma_0^2 = 1/2$ and take $\delta > 0$ to a numerically small value. For the other terms of equation 4.1, we use the formula giving the KL divergence between multivariate normals. Let \mathcal{N}_0 and \mathcal{N}_1 be ℓ -variate normal distributions with mean μ_0, μ_1 and covariance Σ_0, Σ_1 , respectively. Then,

$$\text{KL}(\mathcal{N}_0, \mathcal{N}_1) = \frac{1}{2} \left\{ \text{Tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) - \ell + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right\}.$$

By using this identity, we find the second term of equation 4.1,

$$\begin{aligned} \text{KL}[q_U(z|x_i), q(z|x_i)] &= \frac{1}{2} \left\{ \frac{m\sigma^2 + (\ell - m)\delta^2}{\gamma^2} + \frac{1}{\gamma^2} \|\phi_\theta(x_i) - \mathbb{P}_U \phi_\theta(x_i)\|_2^2 \right. \\ &\quad \left. - \ell + \log \left(\frac{\gamma^{2\ell}}{\sigma^{2m} \delta^{2(\ell-m)}} \right) \right\}, \end{aligned}$$

Table 2: Data Sets and Hyperparameters Used for the Experiments.

Data Set	N	d	m	M
MNIST	60,000	28×28	10	256
fMNIST	60,000	28×28	10	256
SVHN	73,257	$32 \times 32 \times 3$	10	256
dSprites	737,280	64×64	5	256
3DShapes	480,000	$64 \times 64 \times 3$	6	256
3D cars	17,664	$64 \times 64 \times 3$	3	256

Note: N is the number of training samples, d the input dimension (resized images), m the subspace dimension, and M the minibatch size.

where $q(z|x) = \mathcal{N}(z|\phi_\theta(x), \gamma^2 \mathbb{I}_\ell)$. For the third term in equation 4.1, we find

$$\text{KL}[q_U(z|x_i), p(z)] = \frac{1}{2} \left\{ \text{Tr}((\sigma^2 \mathbb{P}_U + \delta^2 \mathbb{P}_{U^\perp}) \Sigma^{-1}) + (\mathbb{P}_U \phi_\theta(x_i))^\top \Sigma^{-1} (\mathbb{P}_U \phi_\theta(x_i)) \right. \\ \left. + \log \det(\Sigma) - \ell - \log(\sigma^{2m} \delta^{2(\ell-m)}) \right\},$$

with $p(z) = \mathcal{N}(0, \Sigma)$. By averaging over $i = 1, \dots, n$, we obtain

$$\frac{1}{n} \sum_{i=1}^n \text{KL}[q_U(z|x_i), p(z)] = \frac{1}{2} \left\{ \text{Tr}((\sigma^2 \mathbb{P}_U + \delta^2 \mathbb{P}_{U^\perp}) \Sigma^{-1}) + \text{Tr}(\mathbb{P}_U C_\theta \mathbb{P}_U \Sigma^{-1}) \right. \\ \left. + \log \det(\Sigma) - \ell - \log(\sigma^{2m} \delta^{2(\ell-m)}) \right\},$$

where we used the cyclic property of the trace and $C_\theta = \frac{1}{n} \sum_{i=1}^n \phi_\theta(x_i) \phi_\theta(x_i)^\top$. This proves the analogous expression in section 4. Finally, the estimation of the optimal Σ can be done in parallel to the maximum likelihood estimation of the covariance matrix of a multivariate normal.

A.3 Data Sets and Hyperparameters. We refer to Tables 2 and 3 for specific details on the model architectures, data sets, and hyperparameters used in this article. All models were trained on full data sets and for a maximum of 1000 epochs. Furthermore, all data sets are scaled between [0-1] and are resized to 28×28 dimensions except dSprites and 3D cars. The PyTorch library (single precision) in Python was used as the programming language on 8 GB NVIDIA QUADRO P4000 GPU. See algorithm 1 for training the St-RKM model. In the case of FactorVAE, the discriminator architecture is same as proposed in the original paper (Kim & Mnih, 2018).

A.3.1 Disentanglement Metrics. MIG was originally proposed by Chen et al. (2018); however, we use the modified metric as proposed in Locatello

Table 3: Model Architectures.

Data Set	Architecture
MNIST/fMNIST/ /SVHN/3DShapes/ sDprites/3DCars	$\phi_{\theta}(\cdot) = \begin{cases} \text{Conv } [c] \times 4 \times 4; \\ \text{Conv } [c \times 2] \times 4 \times 4; \\ \text{Conv } [c \times 4] \times \hat{k} \times \hat{k}; \\ \text{FC 256}; \\ \text{FC 50 (Linear)} \end{cases} \quad \psi_{\xi}(\cdot) = \begin{cases} \text{FC 256}; \\ \text{FC } [c \times 4] \times \hat{k} \times \hat{k}; \\ \text{Conv } [c \times 2] \times 4 \times 4; \\ \text{Conv } [c] \times 4 \times 4; \\ \text{Conv } [c] \text{ (Sigmoid)} \end{cases}$

Notes: All convolutions and transposed convolutions are with stride 2 and padding 1. Unless stated otherwise, layers have parametric-RELU ($\alpha = 0.2$) activation functions, except output layers of the preimage maps, which have sigmoid activation functions (since input data are normalized $[0, 1]$). Adam and Cayley ADAM optimizers have learning rates 2×10^{-4} and 10^{-4} , respectively. The preimage map/decoder network is always taken as transposed of the feature map/encoder network. $c = 48$ for 3D cars; and $c = 64$ for all others. Further, $\hat{k} = 3$ and stride 1 for MNIST, fMNIST, SVHN and 3DShapes; and $\hat{k} = 4$ for others. SVHN and 3DShapes are resized to 28×28 input dimensions.

et al. (2019). We evaluate this score on 5000 test points across all the considered data sets. SAP and Eastwood’s metrics use different classifiers to compute the importance of each dimension of the learned representation for predicting a ground-truth factor. For these metrics, we randomly sample 5000 and 3000 training and testing points, respectively. To compute these metrics, we use the open source library available at github.com/google-research/disentanglement_lib.

A.4 Ablation Studies.

A.4.1 Significance of the KPCA Loss. In this section, we show an ablation study on the KPCA loss and evaluate its effect on disentanglement. We repeat the experiments of section 5 on the mini-3DShapes data set (floor hue, wall hue, object hue, and scale: 8000 samples), where we consider three different variants of the proposed model:

1. **St-RKM** ($\sigma = 0$): The KPCA loss is optimized in a stochastic manner using the Cayley ADAM optimizer, as proposed in this article.
2. **Gen-RKM**: The KPCA loss is optimized exactly at each step by performing an eigendecomposition in each minibatch (this corresponds to the algorithm in Pandey et al., 2021).
3. **AE-PCA**: A standard AE is used, and a reconstruction loss is minimized for the training. As a postprocessing step, a PCA is performed on the latent embedding of the training data.

The encoder/decoder maps are the same across all the models, and for the AE-PCA model, additional linear layers are used to map the latent space to

Table 4: Training Timings per Epoch (in minutes) and Disentanglement Scores (Heusel et al., 2017) for Different Variants of RKM When Trained on the *mini-3Dshapes* Data Set.

		St-RKM ($\sigma = 0$)	Gen-RKM	AE-PCA
Training time		3.01 (0.71)	9.21 (0.54)	2.87 (0.33)
Disentanglement score	Lasso	0.40 (0.02)	0.44 (0.01)	0.35 (0.01)
	RF	0.27 (0.01)	0.31 (0.02)	0.22 (0.02)
Compliance score	Lasso	0.64 (0.01)	0.51 (0.01)	0.42 (0.01)
	RF	0.67 (0.02)	0.58 (0.01)	0.45 (0.02)
Information score	Lasso	1.01 (0.02)	1.11 (0.02)	1.20 (0.01)
	RF	0.98 (0.01)	1.09 (0.01)	1.17 (0.02)

Notes: Gen-RKM has the worst training time but gets the highest disentanglement scores. This is due to the exact eigendecomposition of the kernel matrix at every iteration. This computationally expensive step is approximated by the St-RKM model, which achieves significant speed-up and scalability to large data sets. Finally, the AE-PCA model has the fastest training time due to the absence of eigendecompositions in the training loop. However, using PCA in the post-processing step alters the basis of the latent space. This basis is unknown to the decoder network, resulting in degraded disentanglement performance.

Table 5: FID Scores Computed on Randomly Generated 8000 Images When Trained with Architecture and Hyperparameters.

	St-RKM	VAE	β -VAE	FactorVAE	InfoGAN
MNIST	24.63 (0.22)	36.11 (1.01)	42.81 (2.01)	35.48 (0.07)	45.74 (2.93)
fMNIST	61.44 (1.02)	73.47 (0.73)	75.21 (1.11)	69.73 (1.54)	84.11 (2.58)

Notes: Lower is better with standard deviations. Adapted from Dupont (2018).

the subspace. From Table 4, we conclude that optimizing the KPCA loss during training improves disentanglement. Moreover, using a stochastic algorithm improves computation time and scalability with only a slight decrease in disentanglement score. Note that calculating the exact eigendecomposition at each step (Gen-RKM) comes with numerical difficulties. In particular, double floating-point precision has to be used together with a careful selection of the number of principal components to avoid ill-conditioned kernel matrices. This problem is not encountered when using the St-RKM training algorithm.

A.4.2 Smaller Encoder/Decoder Architecture. In this section, we analyze the impact of the encoder/decoder architecture on the generation quality of considered models. The generation quality experiment of section 5 is repeated on the fMNIST and MNIST data set, where the architecture and hyperparameters are adapted from Dupont (2018). From Table 5 and Figure 9,

Table 6: Computing the Diagonalization Scores (see Figure 3).

Models	dSprites	3DShapes	3D cars
St-RKM-sl ($\sigma = 10^{-3}, U_*$)	0.17 (0.05)	0.23 (0.03)	0.21 (0.04)
St-RKM ($\sigma = 10^{-3}, U_*$)	0.26 (0.05)	0.30 (0.10)	0.31 (0.09)
St-RKM ($\sigma = 10^{-3}, \text{random } U$)	0.61 (0.02)	0.72 (0.01)	0.69 (0.03)

Notes: Denote $M = \frac{1}{|C|} \sum_{i \in C} U_*^\top \nabla \psi(y_i) \nabla \psi(y_i)^\top U_*$, with $y_i = \mathbb{P}_U \phi_\theta(x_i)$ (cf. equation 3.6). Then we compute the score as $\|M - \text{diag}(M)\|_F / \|M\|_F$, where $\text{diag} : \mathbb{R}^{m \times m} \mapsto \mathbb{R}^{m \times m}$ sets the off-diagonal elements of matrix to zero. The scores are computed for each model over 10 random seeds and show the mean (standard deviation). Lower scores indicate better diagonalization.

we see that the overall FID scores and generation quality have improved; however, the relative scores among the models did not change significantly.

A.4.3 Analysis of St-RKM with a Fixed U. We discuss here the role of the optimization of $\text{St}(\ell, m)$ on disentanglement in the case of a classical AE loss ($\sigma = 0$). To do so, a matrix $\tilde{U} \in \text{St}(\ell, m)$ is generated randomly³ and kept fixed during the training of the following optimization problem,

$$\min_{\theta, \xi} \lambda \frac{1}{n} \sum_{i=1}^n L_{\xi, \tilde{U}}^{(0)}(x_i, \phi_\theta(x_i)) + \underbrace{\frac{1}{n} \sum_{i=1}^n \|\mathbb{P}_{\tilde{U}^\perp}^{(\varepsilon)} \phi_\theta(x_i)\|_2^2}_{\text{regularized PCA objective}} \tag{A.2}$$

with $\lambda = 1$ and where $\varepsilon \geq 0$ is a regularization constant and where the regularized (or mollified) projector $\mathbb{P}_{\tilde{U}^\perp}^{(\varepsilon)} = \varepsilon(\tilde{U}\tilde{U}^\top + \varepsilon\mathbb{I}_\ell)^{-1}$ is used in order to prevent numerical instabilities. Indeed, if $\varepsilon = 0$, the second term in equation A.2 (PCA term) is not strictly convex as a function of ϕ_θ , since this quadratic form has flat directions along the column subspace of \tilde{U} . Our numerical simulations in single-precision PyTorch with $\varepsilon = 0$ exhibit instabilities, that is, the PCA term in equation A.2 takes negative values during the training. Hence, the regularized projector is introduced so that the PCA quadratic is strongly convex for $\varepsilon > 0$. This instability is not observed in the training of equation 3.3 where U is not fixed. This is one asset of our training procedure using optimization over Stiefel manifold. Explicitly, the regularized projector satisfies the following properties:

- $\mathbb{P}_{\tilde{U}^\perp}^{(\varepsilon)} u_\perp = u_\perp$ for all $u_\perp \in (\text{range}(U))^\perp$,
- $\mathbb{P}_{\tilde{U}^\perp}^{(\varepsilon)} u = \varepsilon u$ for all $u \in \text{range}(U)$.

³Using a random $\tilde{U} \in \text{St}(\ell, m)$ can be interpreted as sketching the encoder map in the spirit of randomized orthogonal systems (ROS) sketches (see Yang, Pilanci, & Wainwright, 2017).

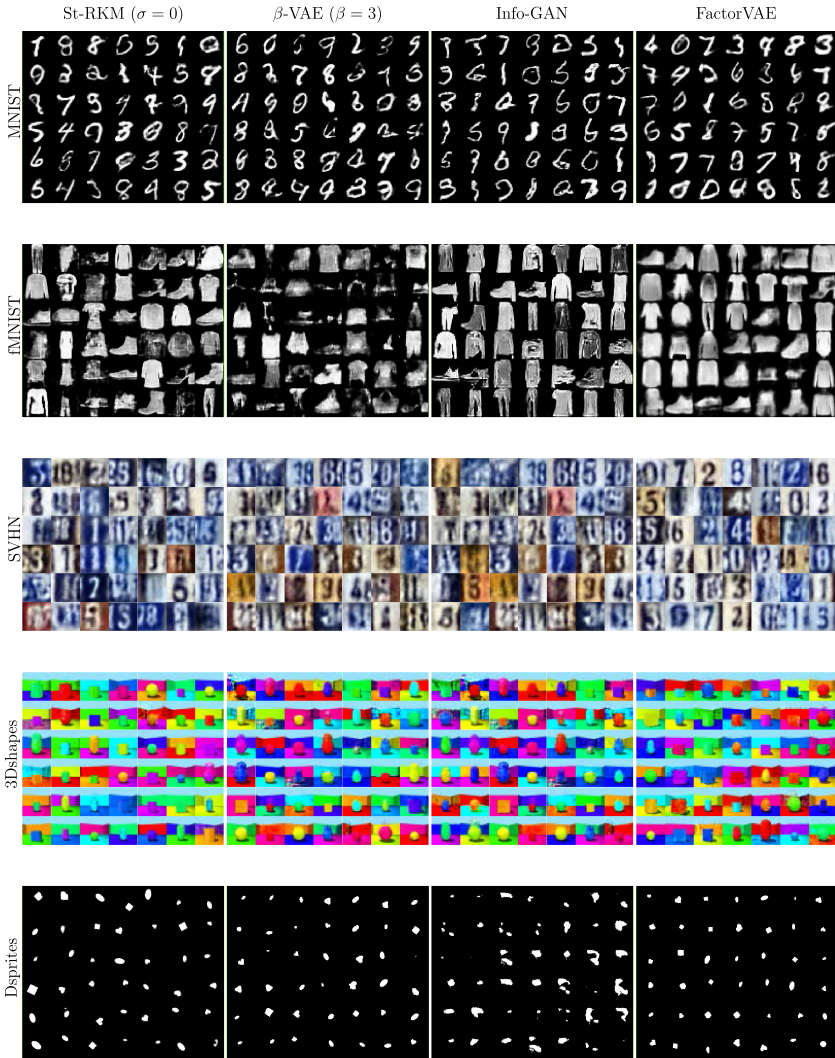


Figure 8: Samples of randomly generated batch of images used to compute FID scores and SWD scores (see Figure 4).

Thanks to the push-through identity, we have the alternative expression $\mathbb{P}_{\hat{U}^\perp}^{(\varepsilon)} = \mathbb{I} - U(U^\top U + \varepsilon \mathbb{I}_m)^{-1}U^\top$. Therefore, it holds $\lim_{\varepsilon \rightarrow 0} \mathbb{P}_{\hat{U}^\perp}^{(\varepsilon)} = \mathbb{P}_{\hat{U}^\perp}$, as it should. In our experiments, we set $\varepsilon = 10^{-5}$. If $\varepsilon \leq 10^{-6}$, the regularized PCA objective in equation A.2 takes negative values after a few epochs due to the numerical instability as mentioned above.

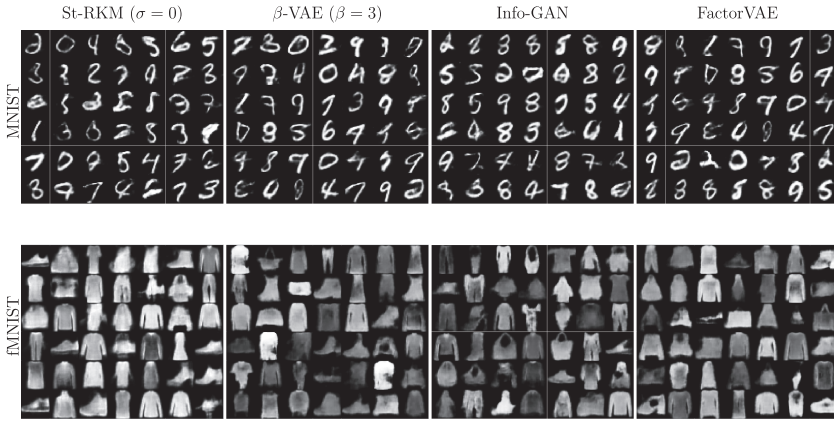
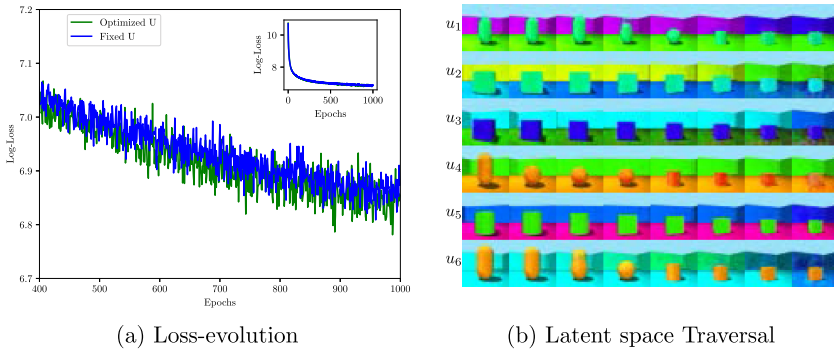


Figure 9: Samples of randomly generated images used to compute the FID scores. See Table 5.



(a) Loss-evolution

(b) Latent space Traversal

Figure 10: (a) Loss evolution (log plot) during the training of equation A.2 over 1000 epochs with $\varepsilon = 10^{-5}$ once with Cayley ADAM optimizer (green curve) and then without (blue curve). (b) Traversals along the principal components when the model was trained with a fixed U , that is, with the objective given by equation A.2 and $\varepsilon = 10^{-5}$. There is no clear isolation of a feature along any of the principal components, indicating further that optimizing over U is key to better disentanglement.

In Figure 10a, the evolution of the training objective A.2 is displayed. It can be seen that the final objective has a lower value [$\exp(6.78) \approx 881$] when U is optimized compared to its fixed counterpart [$\exp(6.81) \approx 905$], showing the merit of optimizing over Stiefel manifold for the same parameter ε . Hence, the subspace determined by $\text{range}(U)$ has to be adapted to the encoder and decoder networks. In other words, the training over θ, ξ is not

sufficient to minimize the $\text{St}(\ell, m)$ objective with Adam. Figure 10b further explores the latent traversals in the context of this ablation study. In the top row of Figure 10b (latent traversal in the direction of \mathbf{u}_1), both the shape of the object and the wall hue are changing. A coupling between wall hue and shape is also visible in the bottom row of this figure.

Acknowledgments

Most of this work was done when M.F. was at KU Leuven.

EU: The research leading to these results received funding from the European Research Council under the European Union's Horizon 2020 research and innovation program/ERC Advanced Grant E-DUALITY (787960). This article reflects only the authors' views, and the EU is not liable for any use that may be made of the contained information.

Research Council KUL: Optimization frameworks for deep kernel machines C14/18/068.

Flemish government: (a) FWO: projects: GOA4917N (Deep Restricted Kernel Machines: Methods and Foundations), PhD/postdoc grant. (b) This research received funding from the Flemish government (AI Research Program). We are affiliated with Leuven.AI-KU Leuven institute for AI, B-3000, Leuven, Belgium.

Ford KU Leuven Research Alliance Project: KUL0076 (stability analysis and performance improvement of deep reinforcement learning algorithms).

Vlaams Supercomputer Centrum: The computational resources and services used in this work were provided by the VSC (Flemish Supercomputer Center), funded by the Research Foundation–Flanders (FWO) and the Flemish government department EWI.

References

- Absil, P.-A., Mahony, R., & Sepulchre, R. (2008). *Optimization algorithms on matrix manifolds*. Princeton, NJ: Princeton University Press.
- Avron, H., Nguyen, H., & Woodruff, D. (2014). Subspace embeddings for the polynomial kernel. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems*, 27 (pp. 2258–2266). Red Hook, NY: Curran.
- Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8), 1798–1828. 10.1109/TPAMI.2013.50, PubMed: 23787338
- Burgess, C., & Kim, H. (2018). *3Dshapes dataset*. <https://github.com/deepmind/3dshapes-dataset/>
- Burgess, C. P., Higgins, I., Pal, A., Matthey, L., Watters, N., Desjardins, G., & Lerchner, A. (2017). Understanding disentangling in β -VAE. In *NIPS 2017 Workshop on Learning Disentangled Representations: From Perception to Control*. <https://sites.google.com/view/disentanglenips2017>

- Chen, R. T. Q., Li, X., Grosse, R. B., & Duvenaud, D. K. (2018). Isolating sources of disentanglement in variational autoencoders. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in neural information processing systems*, 31 (pp. 2610–2620). Red Hook, NY: Curran.
- Dupont, E. (2018). Learning disentangled joint continuous and discrete representations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in neural information processing systems*, 31 (pp. 708–718). Red Hook, NY: Curran.
- Eastwood, C., & Williams, C. K. I. (2018). A framework for the quantitative evaluation of disentangled representations. In *Proceedings of the International Conference on Learning Representations*.
- Ghosh, P., Sajjadi, M. S., Vergari, A., Black, M., & Schölkopf, B. (2020). From variational to deterministic autoencoders. In *Proceedings of the International Conference on Learning Representations*.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., & Hochreiter, S. (2017). GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In I. Guyon, Y. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems*, 30 (pp. 6629–6640). Red Hook, NY: Curran.
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., . . . Lerchner, A. (2017). Beta-VAE: Learning basic visual concepts with a constrained variational framework. In *Proceedings of the International Conference on Learning Representations* (vol. 2, p. 6).
- Hinton, G. E. (2005). What kind of a graphical model is the brain? In *Proceedings of the 19th International Joint Conference on Artificial Intelligence* (pp. 1765–1775).
- Karras, T., Aila, T., Laine, S., & Lehtinen, J. (2017). Progressive growing of GANs for improved quality, stability, and variation. In *Proceedings of the International Conference on Learning Representations*.
- Kim, H., & Mnih, A. (2018). Disentangling by factorising. In *Proceedings of the Thirty-Fifth International Conference on Machine Learning* (vol. 80, pp. 2649–2658).
- Kingma, D. P., & Welling, M. (2014). Auto-encoding variational Bayes. In *Proceedings of the International Conference on Learning Representations*.
- Kumar, A., Sattigeri, P., & Balakrishnan, A. (2018). Variational inference of disentangled latent concepts from unlabeled observations. In *Proceedings of the International Conference on Learning Representations*. <https://openreview.net/forum?id=H1kG7GZAW>
- LeCun, Y., & Cortes, C. (2010). *MNIST handwritten digit database*. <http://yann.lecun.com/exdb/mnist/>
- LeCun, Y., Huang, F. J., & Bottou, L. (2004). Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*.
- Li, J., Li, F., & Todorovic, S. (2020). Efficient Riemannian optimization on the Stiefel manifold via the Cayley transform. In *Proceedings of the International Conference on Learning Representations*.
- Locatello, F., Bauer, S., Lučić, M., Rätsch, G., Gelly, S., Schölkopf, B., & Bachem, O. F. (2019). Challenging common assumptions in the unsupervised learning of

- disentangled representations. In *Proceedings of the International Conference on Machine Learning*.
- Locatello, F., Tschannen, M., Bauer, S., Rätsch, G., Schölkopf, B., & Bachem, O. (2020). Disentangling factors of variations using few labels. In *International Conference on Learning Representations*.
- Matthey, L., Higgins, I., Hassabis, D., & Lerchner, A. (2017). *dSprites: Disentanglement testing Sprites dataset*. <https://github.com/deepmind/dsprites-dataset/>
- Nesterov, Y. (2014). *Introductory lectures on convex optimization: A basic course*. Berlin: Springer.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., & Ng, A. Y. (2011). Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*. http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf
- Pandey, A., Schreurs, J., & Suykens, J. A. K. (2020). Robust generative restricted kernel machines using weighted conjugate feature duality. In *Proceedings of the Sixth International Conference on Machine Learning, Optimization, and Data Science*.
- Pandey, A., Schreurs, J., & Suykens, J. A. (2021). Generative restricted kernel machines: A framework for multi-view generation and disentangled feature learning. *Neural Networks*, 135, 177–191. 10.1016/j.neunet.2020.12.010, PubMed: 33395588
- Reed, S., Zhang, Y., Zhang, Y., & Lee, H. (2015). Deep visual analogy-making. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems*, 28. Red Hook, NY: Curran.
- Rezende, D. J., & Mohamed, S. (2015). Variational inference with normalizing flows. In *Proceedings of the International Conference on Machine Learning*.
- Rolínek, M., Zietlow, D., & Martius, G. (2019). Variational autoencoders pursue PCA directions (by accident). In *Proceedings of the 2019 IEEE/CVF conference on Computer Vision and Pattern Recognition* (pp. 12398–12407).
- Salakhutdinov, R., & Hinton, G. (2009). Deep Boltzmann machines. In *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*.
- Suykens, J. A. K. (2017). Deep restricted kernel machines using conjugate feature duality. *Neural Computation*, 29(8), 2123–2163. 10.1162/neco_a_00984, PubMed: 28562217
- Xiao, H., Rasul, K., & Vollgraf, R. (2017). *Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms*. arXiv:1708.07747.
- Yang, Y., Pilanci, M., & Wainwright, M. J. (2017). Randomized sketches for kernels: Fast and optimal nonparametric regression. *Annals of Statistics*, 45(3), 991–1023. 10.1214/16-AOS1472

Received October 4, 2021; accepted May 24, 2022.