

Variational Autoencoders Pursue PCA Directions (by Accident)

Michal Rolínek*, Dominik Zietlow* and Georg Martius
 Max-Planck-Institute for Intelligent Systems, Tübingen, Germany
 {mrolinek, dzietlow, gmartius}@tue.mpg.de

Abstract

The Variational Autoencoder (VAE) is a powerful architecture capable of representation learning and generative modeling. When it comes to learning interpretable (disentangled) representations, VAE and its variants show unparalleled performance. However, the reasons for this are unclear, since a very particular alignment of the latent embedding is needed but the design of the VAE does not encourage it in any explicit way. We address this matter and offer the following explanation: the diagonal approximation in the encoder together with the inherent stochasticity force local orthogonality of the decoder. The local behavior of promoting both reconstruction and orthogonality matches closely how the PCA embedding is chosen. Alongside providing an intuitive understanding, we justify the statement with full theoretical analysis as well as with experiments.

1. Introduction

The Variational Autoencoder (VAE) [24, 36] is one of the foundational architectures in modern-day deep learning. It serves both as a generative model as well as a representation learning technique. The generative model is predominantly exploited in computer vision [25, 15, 22, 16] with notable exceptions such as generating combinatorial graphs [26]. As for representation learning, there is a variety of applications, ranging over image interpolation [19], one-shot generalization [35], language models [43], speech transformation [3], and more. Aside from direct applications, VAEs embody the success of variational methods in deep learning and have inspired a wide range of ongoing research [23, 44].

Recently, unsupervised learning of interpretable latent representations has received a lot of attention. Interpretability of the latent code is an intuitively clear concept. For instance, when representing faces one latent variable would solely correspond to the gender of the person, another to skin tone, yet another to hair color and so forth. Once such

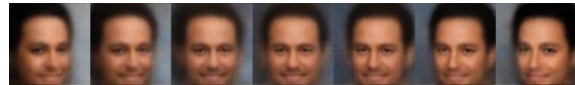


Figure 1. Latent traversal over a single latent coordinate on an exemplary image from the CelebA dataset [28] for a trained β -VAE. The latent coordinate clearly isolates the azimuth angle. Provided by courtesy of the authors of [17].

a representation is found it allows for interpretable latent code manipulation, which is desirable in a variety of applications; recently, for example, in reinforcement learning [39, 18, 11, 41, 34].

The term *disentanglement* [10, 2, 29] offers a more formal approach. A representation is considered disentangled if each latent component encodes precisely one “aspect” (a generative factor) of the data. Under the current disentanglement metrics [17, 21, 6, 29], VAE-based architectures (β -VAE [17], TCVAE [6], FactorVAE [21]) dominate the benchmarks, leaving behind other approaches such as InfoGAN [7] and DCIGN [25]. Exemplarily, a latent traversal for a β -VAE is shown in Fig. 1 in which precisely one generative factor is isolated (face azimuth).

The success of VAE-based architectures on disentanglement tasks comes with a certain surprise. One surprising aspect is that VAEs have been challenged on both of its own design functionalities, as generative models [14, 12] and as log-likelihood optimizers [30, 33]. Yet, no such claims are made in terms of disentanglement. Another surprise stems from the fact that disentanglement requires the following feature: the representative low-dimensional manifold must be aligned well with the coordinate axes. However, the design of the VAE does not suggest any such mechanism. On the contrary, the idealized log-likelihood objective is, for example, invariant to rotational changes in the alignment.

Such observations have planted a suspicion that the inner workings of the VAE are not sufficiently understood. Several recent works approached this issue [5, 40, 8, 1, 12, 31, 9]. However, a mechanistic explanation for the VAE’s unexpected ability to disentangle is still missing.

In this paper, we isolate an internal mechanism of the VAE (also β -VAE) responsible for choosing a particular latent representation and its alignment. We give theoretical

*These authors contributed equally to this work.

analysis covering also the nonlinear case and explain the discovered dynamics intuitively. We show that this mechanism promotes local orthogonality of the embedding transformation and clarify how this orthogonality corresponds to good disentanglement. Further, we uncover strong resemblance between this mechanism and the classical Principle Components Analysis (PCA) algorithm. We confirm our theoretical findings in experiments.

Our theoretical approach is particular in the following ways: (a) we base the analysis on the *implemented* loss function in contrast to the typically considered idealized loss, and (b) we identify a specific regime, prevalent in practice, and utilize it for a vital simplification. This simplification is the crucial step in enabling formalization.

The results, other than being significant on their own, also provide a solid explanation of “why β -VAEs disentangle”.

2. Background

Let us begin with reviewing the basics of VAE, PCA, and of the Singular Value Decomposition (SVD), along with a more detailed overview of disentanglement.

2.1. Variational Autoencoders

Let $\{\mathbf{x}^i\}_{i=1}^N$ be a dataset consisting of N i.i.d. samples $\mathbf{x}^i \in X = \mathbb{R}^n$ of a random variable \mathbf{x} . An autoencoder framework operates with two mappings, the encoder $\text{Enc}_\varphi: X \rightarrow Z$ and the decoder $\text{Dec}_\theta: Z \rightarrow X$, where $Z = \mathbb{R}^d$ is called the *latent space*. In case of the VAE, both mappings are probabilistic and a fixed *prior distribution* $p(\mathbf{z})$ over Z is assumed. Since the distribution of \mathbf{x} is also fixed (actual data distribution $q(\mathbf{x})$), the mappings Enc_φ and Dec_θ induce joint distributions $q(\mathbf{x}, \mathbf{z}) = q_\varphi(\mathbf{z}|\mathbf{x})q(\mathbf{x})$ and $p(x, z) = p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})$, respectively (omitting the dependencies on parameters θ and φ). The idealized VAE objective is then the marginalized log-likelihood

$$\sum_{i=1}^N \log p(\mathbf{x}^i). \quad (1)$$

This objective is, however, not tractable and is approximated by the evidence lower bound (ELBO) [24]. For a fixed \mathbf{x}^i the log-likelihood $\log p(\mathbf{x}^i)$ is lower bounded by

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}^i)} \log p(\mathbf{x}^i | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}^i) \| p(\mathbf{z})), \quad (2)$$

where the first term corresponds to the reconstruction loss and the second to the KL divergence between the latent representation $q(\mathbf{z} | \mathbf{x}^i)$ and the prior distribution $p(\mathbf{z})$. A variant, the β -VAE [17], introduces a weighting β on the KL term for regulating the trade-off between reconstruction (first term) and the proximity to the prior. Our analysis will automatically cover this case as well.

Finally, the prior $p(\mathbf{z})$ is set to $\mathcal{N}(0, \mathcal{I})$ and the encoder is assumed to have the form

$$\text{Enc}_\varphi(\mathbf{x}) \sim q_\varphi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu_\varphi(\mathbf{x}), \text{diag} \sigma_\varphi^2(\mathbf{x})), \quad (3)$$

where μ_φ and σ_φ are deterministic mappings depending on parameters φ . Note particularly, that **the covariance matrix is enforced to be diagonal**. This turns out to be highly significant for the main result of this work. The KL-divergence in (2) can be computed in closed form as

$$L_{\text{KL}} = \frac{1}{2} \sum_{j=1}^d (\mu_j^2(\mathbf{x}^i) + \sigma_j^2(\mathbf{x}^i) - \log \sigma_j^2(\mathbf{x}^i) - 1). \quad (4)$$

In practical implementations, the reconstruction term from (2) is approximated with either a square loss or a cross-entropy loss.

2.2. Disentanglement

In the context of learning interpretable representations [2, 17, 5, 40, 38] it is useful to assume that the data originates from a process with some generating factors. For instance, for images of faces this could be face azimuth, skin brightness, hair length, and so on. Disentangled representations can then be defined as ones in which individual latent variables are sensitive to changes in individual generating factors, while being relatively insensitive to other changes [2]. Although quantifying disentanglement is non-trivial, several metrics have been proposed [21, 17, 6].

Note also, that disentanglement is impossible without first learning a sufficiently expressive latent representation capable of good reconstruction.

In an unsupervised setting, the generating factors are of course unknown and the learning has to resort to statistical properties. Linear dimensionality reduction techniques demonstrate the two basic statistical approaches. Principle Components Analysis (PCA) greedily isolates sources of variance in the data, while Independent Component Analysis (ICA) recovers a factorized representation, see [37] for a recent review.

One important point to make is that **disentanglement is sensitive to rotations of the latent embedding**. Following the example above, let us denote by a , s , and h , continuous values corresponding to face azimuth, skin brightness, and hair length. Then, if we change the ideal latent representation as follows

$$\begin{pmatrix} a \\ s \\ h \end{pmatrix} \mapsto \begin{pmatrix} 0.75a + 0.25s + 0.61h \\ 0.25a + 0.75s - 0.61h \\ -0.61a + 0.61s + 0.50h \end{pmatrix}, \quad (5)$$

we obtain a representation that is equally expressive in terms of reconstruction (in fact we only multiplied with a 3D rotation matrix) but individual latent variables entirely lost their interpretable meaning.

2.3. PCA and Latent Representations

Let us examine more closely how PCA chooses the alignment of the latent embedding and why it matters.

It is well known [4] that for a linear autoencoder with encoder $Y' \in \mathbb{R}^{d \times n}$, decoder $Y \in \mathbb{R}^{n \times d}$, and square error as reconstruction loss, the objective

$$\min_{Y, Y'} \sum_{\mathbf{x}^i \in X} \|\mathbf{x}^i - YY'\mathbf{x}^i\|^2 \quad (6)$$

is minimized by the PCA decomposition. Specifically, by setting $Y' = P_d$, and $Y = P_d^\top$, for $P_d = \mathcal{I}_{d \times n} P \in \mathbb{R}^{d \times n}$, where $P \in \mathbb{R}^{n \times n}$ is an orthogonal matrix formed by the n normalized eigenvectors (ordered by the magnitudes of the corresponding eigenvalues) of the sample covariance matrix of X and $\mathcal{I}_{d \times n} \in \mathbb{R}^{d \times n}$ is a trivial projection matrix.

However, there are many minimizers of (6) that do not induce the same latent representation. In fact, it suffices to append Y' with some invertible transformations (e.g. rotations and scaling) and prefix Y with their inverses. This geometrical intuition is well captured using the singular value decomposition (SVD), see also Figure 2.

Theorem 1 (SVD rephrased, [13]). *Let $M: \mathbb{R}^n \rightarrow \mathbb{R}^d$ be a linear transformation (matrix). Then there exist*

- $U: \mathbb{R}^n \rightarrow \mathbb{R}^n$, an orthogonal transformation (matrix) of the input space,
- $\Sigma: \mathbb{R}^n \rightarrow \mathbb{R}^d$ a “scale-and-embed” transformation (induced by a diagonal matrix),
- $V: \mathbb{R}^d \rightarrow \mathbb{R}^d$, an orthogonal transformation (matrix) of the output space

such that $M = V\Sigma U^\top$.

Remark 1. *For the sake of brevity, we will refer to orthogonal transformations (with slight abuse of terminology) simply as rotations.*

Example 1 (Other minimizers of the PCA objective). *Define Y and Y' with their SVDs as $Y = P^\top \Sigma Q$ and its pseudoinverse $Y' = Y^\dagger = Q^\top \Sigma^\dagger P$ and see that*

$$YY' = P^\top \Sigma Q Q^\top \Sigma^\dagger P = P^\top \mathcal{I}_{d \times n} \mathcal{I}_{n \times d} P = P_d^\top P_d \quad (7)$$

so they are indeed also minimizers of the objective (6) irrespective of our choice of Q and Σ .

It is also straightforward to check that the only choices of Q , which respect the coordinate axes given by PCA, are for $|Q|$ to be a permutation matrix.

The take-away message (valid also in the non-linear case) from this example is:

Different rotations of the same latent space are equally suitable for reconstruction.

Following the PCA example, we formalize which linear mappings have the desired “axes-preserving” property.

Proposition 1 (Axes-preserving linear mappings). *Assume $M \in \mathbb{R}^{n \times d}$ with $d < n$ has d distinct nonzero singular values. Then the following statements are equivalent:*

- The columns of M are (pairwise) orthogonal.*
- In every SVD of M as $M = U\Sigma V^\top$, $|V|$ is a permutation matrix.*

We strongly suggest developing a geometrical understanding for both cases (a) and (b) via Figure 2. For an intuitive understanding of the formal requirement of distinct eigenvalues, we refer to Supp. C.2.

Take into consideration that once the encoder preserves the principle directions of the data, this already ensures an axis-aligned embedding. The same is true also if the decoder is axes-preserving, provided the reconstruction of the autoencoder is accurate.

2.4. Related work

Due to high activity surrounding VAEs, additional care is needed when it comes to evaluating novelty. To the best of our knowledge, two recent works address related questions and require special attention.

The authors of [5] also aim to explain good performance of (β -)VAE in disentanglement tasks. A compelling intuitive picture of the underlying dynamics is drawn and supporting empirical evidence is given. In particular, the authors *hypothesize* that “ β -VAE finds latent components which make different contributions to the log-likelihood term of the cost function [reconstruction loss]”, while suspecting that the diagonal posterior approximation is responsible for this behavior. Our theoretical analysis confirms both conjectures (see Section 4).

Concurrent work [40] develops ISA-VAE; another VAE-based architecture suited for disentanglement. Some parts of the motivation overlap with the content of our work. First, rotationally nonsymmetric priors are introduced for reasons similar to the content of Section 3.1. And second, both orthogonalization and alignment with PCA directions are empirically observed for VAEs applied to toy tasks.

3. Results

3.1. The problem with log-likelihood

The message from Example 1 and from the discussion about disentanglement is clear: latent space *rotation* matters. Let us look how the idealized objectives (1) and (2) handle this.

For a fixed rotation matrix U we will be comparing a baseline encoder-decoder pair ($\text{Enc}_\varphi, \text{Dec}_\theta$) with a pair

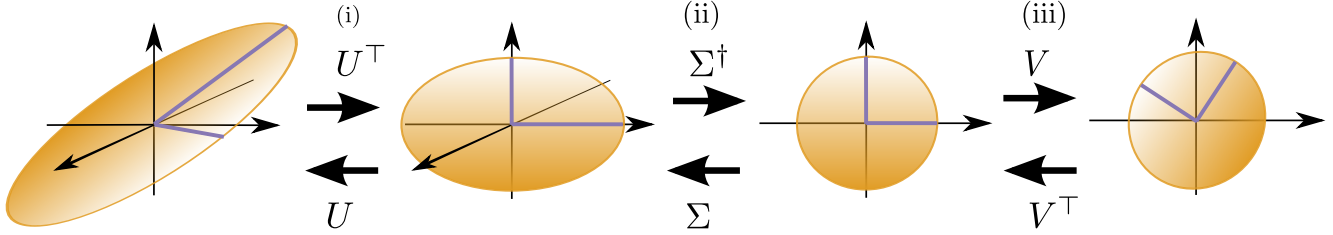


Figure 2. Geometric interpretation of the singular value decomposition (SVD). Sequential illustration of the effects of applying the corresponding SVD matrices of the encoder transformation $V\Sigma^\dagger U^\top$ (left to right) and decoder $U\Sigma V^\top$ (right to left). We notice that steps (i) and (ii) of the encoder preserve the principle directions of the data. Step (iii), however, causes misalignment. In that regard, good encoders are the ones for which step (iii) is trivial. The same argument works for the decoder (in reverse order). This condition is equivalent (for non-degenerate transformations) to $U\Sigma V^\top$ having orthogonal columns (See Proposition 1, where this is phrased for the decoder).

($\text{Enc}_{\varphi,U}, \text{Dec}_{\theta,U}$) defined as

$$\text{Enc}_{\varphi,U}(\mathbf{x}) = U \text{Enc}_{\varphi}(\mathbf{x}), \quad (8)$$

$$\text{Dec}_{\theta,U}(\mathbf{z}) = \text{Dec}_{\theta}(U^\top \mathbf{z}). \quad (9)$$

The shortcomings of idealized losses are summarized in the following propositions.

Proposition 2 (Log-likelihood rotation invariance). *Let φ, θ be any choice of parameters for encoder-decoder pair ($\text{Enc}_{\varphi,U}, \text{Dec}_{\theta,U}$). Then, if the prior $p(\mathbf{z})$ is rotationally symmetric, the value of the log-likelihood objective (1) does not depend on the choice of U .*

Note that the standard prior $\mathcal{N}(0, \mathcal{I})$ is rotationally symmetric. This deficiency is not salvaged by the ELBO approximation.

Proposition 3 (ELBO rotation invariance). *Let φ, θ be any choice of parameters for encoder-decoder pair ($\text{Enc}_{\varphi,U}, \text{Dec}_{\theta,U}$). Then, if the prior $p(\mathbf{z})$ is rotationally symmetric, the value of the ELBO objective (2) does not depend on the choice of U .*

We do not claim novelty of these propositions, however we are not aware of their formalization in the literature. The proofs can be found in Supplementary Material (Suppl. A). An important point now follows:

Log-likelihood based methods (with rotationally symmetric priors) cannot claim to be designed to produce disentangled representations.

However, **enforcing a diagonal posterior of the VAE encoder (3) disrupts the rotational symmetry** and consequently the resulting objective (4) escapes the invariance arguments. Moreover, as we are about to see, this diagonalization comes with beneficial effects regarding disentanglement. We assume this diagonalization was primarily introduced for different reasons (tractability, computational convenience), hence the “by accident” part of the title.

3.2. Reformulating VAE loss

The fact that VAEs were *not meant* to promote orthogonality reflects in some technical challenges. For one, we cannot follow a usual workflow of a theoretical argument; set up an idealized objective and find suitable approximations which allow for stochastic gradient descent (a top-down approach). We need to do the exact opposite, start with the *implemented loss function* and find the right simplifications that allow isolating the effects in question while preserving the original training dynamics (a bottom-up approach). This is the main content of this section.

First, we formalize the typical situation in which VAE architectures “shut down” (fill with pure noise) a subset of latent variables and put high precision on the others.

Definition 1. *We say that parameters φ, θ induce a polarized regime if the latent coordinates $\{1, 2, \dots, d\}$ can be partitioned as $V_a \cup V_p$ (sets of active and passive variables) such that*

$$(a) \mu_j^2(\mathbf{x}) \ll 1 \text{ and } \sigma_j^2(\mathbf{x}) \approx 1 \text{ for } j \in V_p,$$

$$(b) \sigma_j^2(\mathbf{x}) \ll 1 \text{ for } j \in V_a,$$

(c) *The decoder ignores the passive latent components, i.e.*

$$\frac{\partial \text{Dec}_{\theta}(z)}{\partial z_j} = 0 \quad \forall j \in V_p.$$

The polarized regime simplifies the loss L_{KL} from (4); part (a) ensures zero loss for passive variables and part (b) implies that $\sigma_j^2(\mathbf{x}) \ll -\log(\sigma_j^2(\mathbf{x}))$. All in all, the per-sample-loss reduces to

$$L_{\approx \text{KL}}(\mathbf{x}^i) = \frac{1}{2} \sum_{j \in V_a} (\mu_j^2(\mathbf{x}^i) - \log(\sigma_j^2(\mathbf{x}^i)) - 1). \quad (10)$$

We will assume the VAE operates in the polarized regime. In Section 5.2, we show on multiple tasks and datasets that the two objectives align very early in the training. This behavior is well-known to practitioners.

Also, we approximate the reconstruction term in (2), as it is most common, with a square loss

$$L_{\text{rec}}(\mathbf{x}^i) = \mathbb{E} \|\text{Dec}_\theta(\text{Enc}_\varphi(\mathbf{x}^i)) - \mathbf{x}^i\|^2 \quad (11)$$

where the expectation is over the stochasticity of the encoder. All in all, the loss we will analyze has the form

$$\sum_{\mathbf{x}^i \in X} L_{\text{rec}}(\mathbf{x}^i) + L_{\approx\text{KL}}(\mathbf{x}^i). \quad (12)$$

Moreover, the reconstruction loss can be further decomposed into two parts; deterministic and stochastic. The former is defined by

$$\bar{L}_{\text{rec}}(\mathbf{x}^i) = \|\text{Dec}_\theta(\mu(\mathbf{x}^i)) - \mathbf{x}^i\|^2 \quad (13)$$

and captures the square loss of the mean encoder. Whereas the stochastic loss

$$\hat{L}_{\text{rec}}(\mathbf{x}^i) = \mathbb{E} \|\text{Dec}_\theta(\mu(\mathbf{x}^i)) - \text{Dec}_\theta(\text{Enc}_\varphi(\mathbf{x}^i))\|^2 \quad (14)$$

is purely induced by the noise injected in the encoder.

Proposition 4. *If the stochastic estimate $\text{Dec}_\theta(\text{Enc}_\varphi(\mathbf{x}^i))$ is unbiased around $\text{Dec}_\theta(\mu(\mathbf{x}^i))$, then*

$$L_{\text{rec}}(\mathbf{x}^i) = \bar{L}_{\text{rec}}(\mathbf{x}^i) + \hat{L}_{\text{rec}}(\mathbf{x}^i). \quad (15)$$

This decomposition resembles the classical bias-variance decomposition of the square error [20].

3.3. The main result

Now, we finally give theoretical evidence for the central claim of the paper:

Optimizing the stochastic part of the reconstruction loss promotes local orthogonality of the decoder.

On that account, we set up an optimization problem which allows us to optimize the stochastic loss (14) independently of the other two. This will isolate its effects on the training dynamics.

In order to make statements about local orthogonality, we introduce for each \mathbf{x}^i the Jacobian (linear approximation) J_i of the decoder at point $\mu(\mathbf{x}^i)$, i.e.

$$J_i = \frac{\partial \text{Dec}_\theta(\mu(\mathbf{x}^i))}{\partial \mu(\mathbf{x}^i)}.$$

Since, according to (3), the encoder can be written as $\text{Enc}_\varphi(\mathbf{x}^i) = \mu(\mathbf{x}^i) + \varepsilon(\mathbf{x}^i)$ with

$$\varepsilon(\mathbf{x}^i) \sim \mathcal{N}(0, \text{diag } \sigma^2(\mathbf{x}^i)), \quad (16)$$

we can approximate the stochastic loss (14) with

$$\begin{aligned} & \mathbb{E}_{\varepsilon(\mathbf{x}^i)} \left\| \text{Dec}_\theta(\mu(\mathbf{x}^i)) - (\text{Dec}_\theta(\mu(\mathbf{x}^i)) + J_i \varepsilon(\mathbf{x}^i)) \right\|^2 \\ &= \mathbb{E}_{\varepsilon(\mathbf{x}^i)} \|J_i \varepsilon(\mathbf{x}^i)\|^2, \end{aligned} \quad (17)$$

Although we aim to fix the deterministic loss (13), we do not need to freeze the mean encoder and the decoder entirely. Following Example 1, for each J_i and its SVD $J_i = U_i \Sigma_i V_i^\top$, we are free to modify V_i as long we correspondingly (locally) modify the mean encoder.

Then we state the optimization problem as follows:

$$\min_{V_i, \sigma_j^i > 0} \sum_{\mathbf{x}^i \in X} \log \mathbb{E}_{\varepsilon(\mathbf{x}^i)} \|J_i \varepsilon(\mathbf{x}^i)\|^2 \quad (18)$$

$$\text{s. t.} \quad \sum_{\mathbf{x}^i \in X} L_{\approx\text{KL}}(\mathbf{x}^i) = C, \quad (19)$$

where $\varepsilon(\mathbf{x}^i)$ are sampled as in (16).

A few remarks are now in place.

- This optimization is not over network parameters, rather directly over the values of all V_i, σ_j^i (only constrained by (19)).
- Both the objective and the constraint concern *global losses*, not per sample losses.
- Indeed, none of V_i, σ_j^i interfere with the rest of the VAE objective (12).

The presence of the (monotone) log function has one main advantage; we can describe **all global minima** of (18) in closed form. This is captured in the following theorem, the technical heart of this work.

Theorem 2 (Main result). *The following holds for optimization problem (18, 19):*

- Every local minimum is a global minimum.
- In every global minimum, the columns of every J_i are orthogonal.

The full proof as well as an explicit description of the minima is given in Suppl. A.1. However, an outline of the main steps is given in the next section on the example of a linear decoder.

The presence of the log term in (18) admittedly makes our argument indirect. There are, however, a couple of points to make. First, as was mentioned earlier, encouraging orthogonality was *not a design feature* of the VAE. In this sense, it is unsurprising that our results are also mildly indirect.

Also, and more importantly, the global optimality of Theorem 2 also implies that, locally, orthogonality is encouraged even for the pure (without logarithm) stochastic loss.

Corollary 1. *For fixed $\mathbf{x}^i \in X$ consider a subproblem of (18) defined as*

$$\min_{V_i, \sigma_j^i > 0} \mathbb{E}_{\varepsilon(\mathbf{x}^i)} \|J_i \varepsilon(\mathbf{x}^i)\|^2 \quad (20)$$

$$\text{s. t.} \quad L_{\approx\text{KL}}(\mathbf{x}^i) = C_i. \quad (21)$$

Also then, the result on the structure of local (global) minima holds:

- (a) Every local minimum is a global minimum.
- (b) In every global minimum, the columns of every J_i are orthogonal.

All in all, Theorem 2 justifies the central message of the paper stated at the beginning of this section. The analogy with PCA is now also clearer. Locally, VAEs optimize a tradeoff between reconstruction and orthogonality.

This result is unaffected by the potential β term in Equation (2), although an appropriate β might be required to ensure the polarized regime.

4. Proof outline

In this section, we sketch the key steps in the proof of Theorem 2 and, more notably, the intuition behind them. The full proof can be found in Suppl. A.1.

We will restrict ourselves to a simplified setting. Consider a linear decoder M with SVD $M = U\Sigma V^T$, which removes the necessity of local linearization. This reduces the objective (18) from a “global” problem over all examples \mathbf{x}^i to an objective where we have the same subproblem for each \mathbf{x}^i .

As in optimization problem (18, 19), we resort to fixing the mean encoder (imagine a well performing one).

In the next paragraphs, we separately perform the optimization over the parameters σ and the optimization over the matrix V .

4.1. Weighting precision

For this part, we fix the decoder matrix M and optimize over values $\sigma^2 = (\sigma_1^2, \dots, \sigma_d^2)$. The simplified objective is

$$\min_{\sigma} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, \text{diag}(\sigma^2))} \|M\varepsilon\|^2 \quad (22)$$

$$\text{s. t.} \quad \sum_j -\log \sigma_j^2 = C, \quad (23)$$

where the $\|\mu\|^2$ terms from (10) disappear since the mean encoder is fixed.

The values $-\log(\sigma_j)$ can now be thought of as precisions allowed for different latent coordinates. The log functions even suggests thinking of the number of significant digits. Problem (22) then asks to distribute the “total precision budget” so that the deviation from decoding “uncorrupted” values is minimal.

We will now solve this problem on an example linear decoder $M_1: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ given by

$$M_1: \begin{pmatrix} x \\ y \end{pmatrix} \mapsto \begin{pmatrix} 4x + y \\ -3x + y \\ 5x - y \end{pmatrix}. \quad (24)$$

Already here we see, that the latent variable x seems more influential for the reconstruction. We would expect that x receives higher precision than y .

Now, for $\varepsilon = (\varepsilon_x, \varepsilon_y)$, we compute

$$\|M_1\varepsilon\|^2 = \|4\varepsilon_x + \varepsilon_y\|^2 + \|-3\varepsilon_x + \varepsilon_y\|^2 + \|5\varepsilon_x - \varepsilon_y\|^2$$

and after taking the expectation, we can use the fact that ε has zero mean and write

$$\mathbb{E} \|M_1\varepsilon\|^2 = \text{var}(4\varepsilon_x + \varepsilon_y) + \text{var}(-3\varepsilon_x + \varepsilon_y) + \text{var}(5\varepsilon_x - \varepsilon_y).$$

Finally, we use that for uncorrelated random variables A and B we have $\text{var}(A + cB) = \text{var} A + c^2 \text{var} B$. After rearranging we obtain

$$\begin{aligned} \mathbb{E} \|M_1\varepsilon\|^2 &= \sigma_x^2(4^2 + (-3)^2 + 5^2) + \sigma_y^2(1^2 + 1^2 + (-1)^2) \\ &= 50\sigma_x^2 + 3\sigma_y^2, \end{aligned}$$

where $\sigma = (\sigma_x^2, \sigma_y^2)$. Note that the coefficients are the **squared norms of the column vectors** of M_1 .

This turns the optimization problem (22) into a simple exercise, particularly after realizing that (23) fixes the value of the product $\sigma_x\sigma_y$. Indeed, we can even set $a^2 = 50\sigma_x$ and $b^2 = 3\sigma_y$ in the trivial inequality $a^2 + b^2 \geq 2ab$ and find that

$$\mathbb{E} \|M_1\varepsilon\|^2 = 50\sigma_x^2 + 3\sigma_y^2 \geq 2 \cdot \sqrt{50 \cdot 3} \cdot e^{-C} \approx 24.5e^{-C}, \quad (25)$$

with equality achieved when $\sigma_x^2/\sigma_y^2 = 3/50$. This also implies that the precision $-\log \sigma_x^2$ on variable x will be considerably higher than for y , just as expected.

Two remarks regarding the general case follow.

- The full version of inequality (25) relies on the concavity of the log function; in particular, on (a version of) Jensen’s inequality.
- The minimum value of the objective depends on the product of the column norms. This also carries over to the unsimplified setting.

4.2. Isolating sources of variance

Now that we can find optimal values of precision, the focus changes on optimally rotating the latent space. In order to understand how such rotations influence the minimum of objective (22), let us consider the following example in which we again resort to decoder matrix $M_2: \mathbb{R}^2 \rightarrow \mathbb{R}^3$.

Imagine, the encoder alters the latent representation by a 45° rotation. Then we can adjust the decoder M_1 by first undoing this rotation. In particular, we set $M_2 = M_1 R_{45^\circ}^\top$, where R_θ is a 2D rotation matrix, rotating by angle θ . We have

$$M_2: \begin{pmatrix} x' \\ y' \end{pmatrix} \mapsto \begin{pmatrix} \frac{1}{2}\sqrt{2}(3x' + 5y') \\ \sqrt{2}(-2x' - y') \\ \sqrt{2}(3x' + 2y') \end{pmatrix}$$

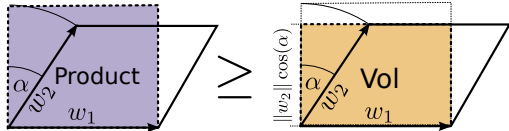


Figure 3. 2D illustration of orthogonality in MV^\top . The vectors w_1, w_2 are the columns of MV^\top . Minimizing the product $\|w_1\| \|w_2\|$ while maintaining the volume $\|w_1\| \|w_2\| \cos(\alpha)$ results in $w_1 \perp w_2$.

and performing analogous optimization as before gives

$$\mathbb{E} \|M_2 \varepsilon\|^2 = \frac{61}{2} \sigma_x^2 + \frac{45}{2} \sigma_y^2 \geq 2 \sqrt{\frac{61 \cdot 45}{4}} e^{-C} \approx 52.4 e^{-C}. \quad (26)$$

We see that the minimal value of the objective is more than twice as high, a substantial difference. On a high level, the reason M_1 was a better choice of a decoder is that the variables x and y had very different impact on the reconstruction. This allowed to save some precision on variable y , as it had smaller effect, and use it on x , where it is more beneficial.

For a higher number of latent variables, one way to achieve a “maximum stretch” among the impacts of latent variables, is to pick them greedily, always picking the next one so that its impact is maximized. This is, at heart, the greedy algorithm for PCA.

Let us consider a slightly more technical statement. We saw in (25) and (26) that after finding optimal values of σ the remaining objective is the product of the column norms of matrix M . Let us denote such quantity by $\text{col}_\Pi(M) = \prod_j \|M_{\cdot j}\|$. Then for a fixed matrix M , we optimize

$$\min_V \text{col}_\Pi(MV^\top) \quad (27)$$

over orthogonal matrices V .

This problem can be interpreted geometrically. The column vectors of MV^\top are the images of base vectors e_j . Consequently, the product gives an upper bound on the volume (the image of the unit cube)

$$\prod_j \|MV^\top e_j\| \geq \text{Vol}(\{MV^\top x : x \in [0, 1]^d\}). \quad (28)$$

However, as orthogonal matrices V are isometries, they do not change this volume. Also, the bound (28) is tight precisely when the vectors $MV^\top e_j$ are orthogonal. Hence, the only way to optimize $\text{col}_\Pi(MV^\top)$ is by tightening the bound, that is by finding V for which the column vectors of MV^\top are orthogonal, see Figure 3 for an illustration. In this regards, it is important that M performs a different scaling along each of the axis (using Σ), which allows for changing the angles among the vectors $MV^\top e_j$ (cf. Figure 2).

Table 1. Percentage of training time where $\Delta_{KL} < 3\%$ (Eq. (30)) continuously until the end. Reported for β -VAE with low (dataset dependent) and high (10) latent dimension.

	β -VAE (dep.)	β -VAE (10)
dSprites	97.8 %	90.6 %
fMNIST	99.8 %	97.7 %
MNIST	99.8 %	99.5 %
Synth. Lin.	99.8 %	96.7 %
Synth. Non-Lin.	99.9 %	98.5 %

5. Experiments

We performed several experiments with different architectures and datasets to validate our results empirically. We show the prevalence of the polarized regime, the strong orthogonal effects of the (β -)VAE, as well as the links to disentanglement.

5.1. Setup

Architectures. We evaluate the classical VAE, β -VAE, a plain autoencoder, and β -VAE $_\Sigma$, where the latter removes the critical diagonal approximation (3) and produces a full covariance matrix $\Sigma(\mathbf{x}^i)$ for every sample. The resulting KL term of the loss is changed accordingly (see Suppl. B.3 for details).

Datasets. We evaluate on the well-known datasets dSprites [32], MNIST [27] and FashionMNIST [42], as well as on two synthetic ones. For both synthetic tasks the input data X is generated by embedding a unit square $V = [0, 1]^2$ into a higher dimension. The latent representation is then expected to be disentangled with respect to axes of V . In one case (*Synth. Lin.*) we used a linear transformation $f_{\text{lin}}: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ and in the other one a non-linear (*Synth. Non-Lin.*) embedding $f_{\text{non-lin}}: \mathbb{R}^2 \rightarrow \mathbb{R}^6$. The exact choice of transformations can be found in Suppl. B. Further information regarding network structures and training parameters is also provided in Suppl. B.4.

Disentanglement metric. For quantifying the disentanglement of a representation, the so called Mutual Information Gap (MIG) was introduced in [6]. As MIG is not well defined for continuous variables, we use an adjusted definition comprising both continuous and discrete variables, simply referred to as *Disentanglement score*. Details are described in Suppl. B.1. Just as in the case of MIG, the Disentanglement score is a number between 0 and 1, where higher value means stronger disentanglement.

Orthogonality metric. For measuring the practical effects of Theorem 2, we introduce a measure of non-orthogonality. As argued in Proposition 1 and Figure 2, for a good decoder M and its SVD $M = U\Sigma V^\top$, the matrix V should be trivial (a signed permutation matrix). We measure the non-triviality with the *Distance to Orthogonality* (DtO) defined as follows. For each \mathbf{x}^i , $i = 1, \dots, N$, employing again the

Table 2. Results for the distance to orthogonality DtO of the decoder (Equation 29) and disentanglement score for different architectures and datasets. Lower DtO values are better and higher Disent. values are better. Random decoders provide a simple baseline for the numbers.

		β -VAE	VAE	AE	β -VAE $_{\Sigma}$	Random Decoder
dSprites	Disent. \uparrow	0.33 \pm 0.15	0.21 \pm 0.10	0.09 \pm 0.04	0.12 \pm 0.06	1.86 \pm 0.11
	DtO \downarrow	0.76 \pm 0.08	1.08 \pm 0.15	1.62 \pm 0.03	1.73 \pm 0.14	
Synth. Lin.	Disent. \uparrow	0.99 \pm 0.01	–	0.71 \pm 0.19	0.71 \pm 0.31	0.79 \pm 0.21
	DtO \downarrow	0.00 \pm 0.00	–	0.33 \pm 0.18	0.34 \pm 0.35	
Synth. Non-Lin.	Disent. \uparrow	0.73 \pm 0.16	–	0.59 \pm 0.30	0.42 \pm 0.24	0.89 \pm 0.16
	DtO \downarrow	0.18 \pm 0.02	–	0.54 \pm 0.13	0.55 \pm 0.02	
MNIST	DtO \downarrow	–	1.59 \pm 0.08	1.83 \pm 0.05	1.93 \pm 0.08	2.11 \pm 0.11
fMNIST	DtO \downarrow	–	1.36 \pm 0.05	1.87 \pm 0.03	2.02 \pm 0.08	2.11 \pm 0.11

Jacobian J_i of the decoder at \mathbf{x}^i and its SVD $J_i = U_i \Sigma_i V_i^T$ and define

$$\text{DtO} = \frac{1}{N} \sum_{i=1}^N \|V_i - P(V_i)\|_F, \quad (29)$$

where $\|\cdot\|_F$ is the Frobenius norm and $P(V_i)$ is a signed permutation matrix that is closest to V (in L^1 sense). Finding the nearest permutation matrix is solved to optimality via mixed-integer linear programming (see Suppl. B.2).

5.2. Polarized regime

In Section 3.2, we assumed VAEs operate in a polarized regime and approximated L_{KL} , the KL term of the implemented objective (4), with $L_{\approx\text{KL}}$ (10). In Table 1 we show that the polarized regime is indeed dominating the training in all examples after a short initial phase. We report the fraction of the training time in which the relative error

$$\Delta_{KL} = \frac{|L_{\text{KL}} - L_{\approx\text{KL}}|}{L_{\text{KL}}} \quad (30)$$

stays below 3% continuously until the end (evaluated every 500 batches). Active variables can be selected by $\sqrt{\text{var}(\mu_j(\mathbf{x}^i))} > 0.5$.

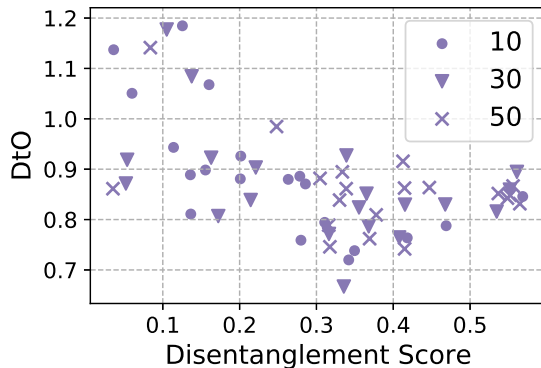


Figure 4. Alignment of the latent representation (low DtO, (29)) results in better disentanglement (higher score). Each datapoint corresponds to an independent run with 10, 30, or 50 epochs.

5.3. Orthogonality and Disentanglement

Now, we provide evidence for Theorem 2 by investigating the DtO (29) for a variety of architectures and datasets, see Table 2. The results clearly support the claim that the VAE based architectures indeed strive for local orthogonality. By generalizing the β -VAE architecture, such that the approximate posterior is any multivariate Gaussian (β -VAE $_{\Sigma}$), the objective becomes rotationally symmetric (just as the idealized objective). As such, no specific alignment is prioritized. The simple autoencoders also do not favor particular orientations of the latent space.

Another important observation is the clear correlation between DtO and the disentanglement score. We show this in Figure 4 where different restarts of the same β -VAE architecture on the dSprites dataset are displayed. We used the state-of-the-art value $\beta = 4$ [17]. Additional experiments are reported in Suppl. C.

6. Discussion

We isolated the mechanism of VAE that leads to local orthogonalization and, in effect, to performing local PCA. Additionally, we demonstrated the functionality of this mechanism in intuitive terms, in formal terms, and also in experiments. We also explained why this behavior is desirable for enforcing disentangled representations.

Our insights show that VAEs make use of the differences in variance to form the representation in the latent space – collapsing to PCA in the linear case. This does not *directly* encourage factorized latent representations. With this in mind, it makes perfect sense that recent improvements of (β -)VAE [6, 21, 40] incorporate additional terms promoting precisely independence.

It is also unsatisfying that VAEs promote orthogonality somewhat indirectly. It would seem that designing architectures allowing explicit control over this feature would be beneficial.

References

- [1] Alexander Alemi, Ben Poole, Ian Fischer, Joshua Dillon, Rif A. Saurous, and Kevin Murphy. Fixing a broken ELBO. In *Proc. 35th Intl. Conference on Machine Learning (ICML)*, volume 80, pages 159–168. PMLR, 2018. [1](#)
- [2] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, Aug 2013. [1](#), [2](#)
- [3] Merlijn Blaauw and Jordi Bonada. Modeling and transforming speech using variational autoencoders. In *INTER-SPEECH*, pages 1770–1774, 2016. [1](#)
- [4] H. Bourlard and Y. Kamp. Auto-association by multilayer perceptrons and singular value decomposition. Manuscript M217, Philips Research Laboratory, Brussels, Belgium, 1987. [3](#)
- [5] Christopher P. Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -vae. *ArXiv e-prints*, abs/1804.03599, 2018. [1](#), [2](#), [3](#)
- [6] Tian Qi Chen, Xuechen Li, Roger B. Grosse, and David K. Duvenaud. Isolating sources of disentanglement in variational autoencoders. *ArXiv e-prints*, abs/1802.04942, 2018. [1](#), [2](#), [7](#), [8](#)
- [7] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel. InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets. *ArXiv e-prints*, June 2016. [1](#)
- [8] B. Dai, Y. Wang, J. Aston, G. Hua, and D. Wipf. Hidden talents of the variational autoencoder. *ArXiv e-prints*, abs/1706.05148, 2018. [1](#)
- [9] Bin Dai and David Wipf. Diagnosing and enhancing vae models. *ArXiv e-prints*, abs/1903.05789, 2019. [1](#)
- [10] Guillaume Desjardins, Aaron Courville, and Yoshua Bengio. Disentangling factors of variation via generative entangling. *ArXiv e-prints*, abs/1210.5474, 2012. [1](#)
- [11] A. Ghadirzadeh, A. Maki, D. Kragic, and M. Bjrkman. Deep predictive policy training using reinforcement learning. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2351–2358, Sept 2017. [1](#)
- [12] Partha Ghosh, Mehdi S. M. Sajjadi, Antonio Vergari, Michael Black, and Bernhard Schölkopf. From variational to deterministic autoencoders. *ArXiv e-prints*, abs/1903.12436, 2019. [1](#)
- [13] G. H. Golub and W. Kahan. Calculating the singular values and pseudo-inverse of a matrix. *Journal of the Society for Industrial and Applied Mathematics: Series B, Numerical Analysis*, 2(2):205–224, 1965. [3](#)
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. [1](#)
- [15] K. Gregor, F. Besse, D. Jimenez Rezende, I. Danihelka, and D. Wierstra. Towards conceptual compression. *ArXiv e-prints*, abs/1604.08772, 2016. [1](#)
- [16] Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. In Francis Bach and David Blei, editors, *Proc. ICML*, volume 37, pages 1462–1471. PMLR, 2015. [1](#)
- [17] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. Beta-VAE: Learning basic visual concepts with a constrained variational framework. *ICLR*, 2017. [1](#), [2](#), [8](#)
- [18] I. Higgins, A. Pal, A. A. Rusu, L. Matthey, C. P Burgess, A. Pritzel, M. Botvinick, C. Blundell, and A. Lerchner. DARLA: Improving zero-shot transfer in reinforcement learning. *ArXiv e-prints*, July 2017. [1](#)
- [19] X. Hou, L. Shen, K. Sun, and G. Qiu. Deep feature consistent variational autoencoder. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1133–1141, March 2017. [1](#)
- [20] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning: With Applications in R*. Springer Publishing Company, Incorporated, 2014. [5](#)
- [21] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas Krause, editors, *Proc. ICML*, volume 80, pages 2649–2658. PMLR, 2018. [1](#), [2](#), [8](#)
- [22] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4743–4751. Curran Associates, Inc., 2016. [1](#)
- [23] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4743–4751. Curran Associates, Inc., 2016. [1](#)
- [24] D. P Kingma and M. Welling. Auto-Encoding Variational Bayes. *ICLR*, 2014. [1](#), [2](#)
- [25] Tejas D Kulkarni, William F. Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2539–2547. Curran Associates, Inc., 2015. [1](#)
- [26] Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. *ArXiv e-prints*, abs/1703.01925, 2017. [1](#)
- [27] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998. [7](#)
- [28] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. [1](#)

- [29] Francesco Locatello, Stefan Bauer, Mario Lucic, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. *ArXiv e-prints*, abs/1811.12359, 2018. [1](#)
- [30] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv e-print*, abs/1511.05644, 2015. [1](#)
- [31] Emile Mathieu, Tom Rainforth, N. Siddharth, and Yee Whye Teh. Disentangling disentanglement in variational auto-encoders. *ArXiv e-prints*, abs/1812.02833, 2018. [1](#)
- [32] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017. [7](#)
- [33] L. Mescheder, S. Nowozin, and A. Geiger. Adversarial variational bayes: Unifying variational autoencoders and generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*. PMLR, Aug. 2017. [1](#)
- [34] A. Nair, V. Pong, M. Dalal, S. Bahl, S. Lin, and S. Levine. Visual Reinforcement Learning with Imagined Goals. *ArXiv e-prints*, abs/1807.04742, July 2018. [1](#)
- [35] Danilo Jimenez Rezende, Shakir Mohamed, Ivo Danihelka, Karol Gregor, and Daan Wierstra. One-shot generalization in deep generative models. *ArXiv e-prints*, abs/1603.05106, 2016. [1](#)
- [36] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *ICML*, 2014. [1](#)
- [37] Karl Ridgeway. A survey of inductive biases for factorial representation-learning. *ArXiv e-prints*, abs/1612.05299, 2016. [2](#)
- [38] Jrgen Schmidhuber. Learning factorial codes by predictability minimization. *Neural Computation*, 4(6):863–879, 1992. [2](#)
- [39] Iulian Vlad Serban, Chinnadhurai Sankar, Mathieu Germain, Saizheng Zhang, Zhouhan Lin, Sandeep Subramanian, Taesup Kim, Michael Pieper, Sarath Chandar, Nan Rosemary Ke, Sai Mudumba, Alexandre de Brébisson, Jose Sotelo, Dendi Suhubdy, Vincent Michalski, Alexandre Nguyen, Joelle Pineau, and Yoshua Bengio. A deep reinforcement learning chatbot. *ArXiv e-prints*, abs/1709.02349, 2017. [1](#)
- [40] Jan Stühmer, Richard Turner, and Sebastian Nowozin. ISA-VAE: Independent subspace analysis with variational autoencoders. In *Submitted to International Conference on Learning Representations*, 2019. [1](#), [2](#), [3](#), [8](#)
- [41] H. van Hoof, N. Chen, M. Karl, P. van der Smagt, and J. Peters. Stable reinforcement learning with autoencoders for tactile and visual data. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3928–3934, Oct 2016. [1](#)
- [42] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *ArXiv e-prints*, abs/1708.07747, 2017. [7](#)
- [43] Zichao Yang, Zhiting Hu, Ruslan Salakhutdinov, and Taylor Berg-Kirkpatrick. Improved variational autoencoders for text modeling using dilated convolutions. In Doina Precup and Yee Whye Teh, editors, *Proc. ICML*, volume 70, pages 3881–3890. PMLR, 2017. [1](#)
- [44] Cheng Zhang, Judith Bütepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. *ArXiv e-prints*, abs/1711.05597, 2017. [1](#)