

# Disentanglement Analysis in Deep Latent Variable Models Matching Aggregate Posterior Distributions

Surojit Saha

SCI, Kahlert School of Computing  
The University of Utah  
Salt Lake City, USA  
surojit.saha@utah.edu

Sarang Joshi

SCI, Kahlert School of Computing  
The University of Utah  
Salt Lake City, USA  
sarang.joshi@utah.edu

Ross Whitaker

SCI, Kahlert School of Computing  
The University of Utah  
Salt Lake City, USA  
whitaker@cs.utah.edu

**Abstract**—Deep latent variable models (DLVMs) are designed to learn meaningful representations in an unsupervised manner, such that the hidden explanatory factors are interpretable by independent latent variables (aka disentanglement). The variational autoencoder (VAE) [1], [2] is a popular DLVM widely studied in disentanglement analysis due to the modeling of the posterior distribution using a factorized Gaussian distribution [3] that encourages the alignment of the latent factors with the latent axes. Several metrics have been proposed recently, assuming that the latent variables explaining the variation in data are aligned with the latent axes (cardinal directions). However, there are other DLVMs, such as the AAE and WAE-MMD (matching the aggregate posterior to the prior), where the latent variables might not be aligned with the latent axes. In this work, we propose a statistical method to *evaluate disentanglement* for any DLVMs in general. The proposed technique discovers the latent vectors representing the generative factors of a dataset that can be different from the cardinal latent axes. We empirically demonstrate the advantage of the method on two datasets.

**Index Terms**—Disentanglement Analysis, Deep Latent Variable Models, Marginal Posterior Matching

## I. INTRODUCTION

Deep latent variable models (DLVMs) have gained a great deal of well-deserved attention due to their ability to *model the distribution* of the high-dimensional, complex datasets [4]–[6] and learn meaningful representations [7]–[12] for downstream applications, such as *few-shot learning* [9], [13], [14]. DLVMs learn a joint distribution,  $p_\theta(\mathbf{x}, \mathbf{z})$ , that captures the relationship between a set of learned, hidden variables,  $\mathbf{z}$ , and the observed variables,  $\mathbf{x}$ . The variational autoencoder (VAE) [1], [2] is a popular DLVM. Learning *disentangled representations* in an unsupervised framework is a desired property of a DLVM such that independent latent variables can explain the variability in the observed data [7]. DLVMs learn disentangled representations by encoding meaningful information onto the independent latent variables, such that each latent variable,  $\mathbf{z}_i$ , represents *only* a single generative factor of the data; thus, making  $\mathbf{z}_i$  *interpretable*.

The DLVM uses an *encoder-decoder* architecture, where the encoder projects the observed data onto a low dimensional manifold, and the decoder reconstructs the encoded representations. The encoded representations are mapped to a prior distribution, and typically, standard normal distribution,  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , is chosen as the prior distribution in DLVMs.

$\mathcal{N}(\mathbf{0}, \mathbf{I})$  is invariant to rotations, implying there is no difference in the expressiveness of a latent sample,  $\mathbf{z}$ , and its rotated counterpart,  $rot(\mathbf{z})$ , in terms of the reconstruction by the decoder. However, the latent representations are not interpretable anymore on rotation [3], [15]. Despite this limitation, the VAE (using  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  as the prior) succeeds in learning a disentangled representation due to the modeling of the posterior distribution in the latent space with a factorized Gaussian distribution [3]. This motivated the development of the different variants of the VAE that encourage the learning of disentangled representations [6], [16]–[18]. In addition, a slew of metrics were proposed to evaluate the disentanglement of the learned representations [16]–[21].

Increasing the strength of the regularization loss in VAEs for improved disentanglement [16] results in the *posterior collapse* [22]–[24] i.e., uninformative latent variables leading to poor reconstruction. The use of additional regularization loss in the VAE objective [17], [18] for better disentanglement often results in the mismatch between the aggregate posterior distribution and the prior [5], [25]. The mismatch leads to the generation of poor-quality samples due to the presence of *pockets/holes* in the encoded distribution. DLVMs other than the VAE, such as the AAE [26], WAE [27], GENs [28], and the AVAE [5] do not suffer from the posterior collapse and the GENs [28], and the methods in [5], [28] closely match the prior. However, matching the aggregate posterior to the prior in these models does not encourage the alignment of the latent generative factors identified by the models with the latent axes, unlike the VAE. Subsequently, resulting in poor performance under the existing disentanglement metrics [16]–[20].

In this paper, we propose a method to *evaluate* disentanglement of any trained DLVM, and not particularly VAEs. The proposed technique identifies directions (unit vectors) in the latent space representing latent variables associated with the true generative factors instead of relying on the cardinal latent axes as in VAEs. Existing metrics for disentanglement analysis use synthetic datasets with *known latent factors* to evaluate the performance of any DLVM [15] as the true latent factors are unknown for real-world datasets. We use the same labels to determine the latent directions (representing the latent factors) for any DLVM. The proposed technique for identifying latent vectors presents a generalized framework that results in better

metric scores across DLVMs, particularly for methods that match aggregate posterior distributions, e.g., the AAE [26], WAE [27], and AVAE [5].

## II. PROPOSED METHOD

### A. Background

The VAE use a probabilistic encoder,  $\mathbf{E}_\phi$ , and a probabilistic decoder,  $\mathbf{D}_\theta$ , to represent  $q_\phi(\mathbf{z} | \mathbf{x})$  and  $p_\theta(\mathbf{x} | \mathbf{z})$ , respectively. Both  $\mathbf{E}_\phi$  and  $\mathbf{D}_\theta$  are usually deep neural networks parameterized by  $\phi$  and  $\theta$ , respectively. The prior distribution,  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  and the surrogate posterior is a factorized Gaussian distribution with diagonal covariance (assuming independent latent dimensions), which is defined as follows:

$$q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\mathbf{x}, \boldsymbol{\sigma}_\mathbf{x}^2 \mathbf{I}), \text{ where } \boldsymbol{\mu}_\mathbf{x}, \boldsymbol{\sigma}_\mathbf{x}^2 \leftarrow \mathbf{E}_\phi(\mathbf{x}). \quad (1)$$

The choice of the Gaussian distribution as the posterior,  $q_\phi(\mathbf{z} | \mathbf{x})$ , helps in efficient computation (reparameterization trick) of the  $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} p_\theta(\mathbf{x} | \mathbf{z})$  in the VAE objective function (ELBO),

$$\max_{\theta, \phi} \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x} | \mathbf{z}) - \text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})) \right]. \quad (2)$$

In contrast, DLVMs matching the aggregate posterior to the prior minimizes the KL ( $q_\phi(\mathbf{z}) || p(\mathbf{z})$ ) in equation 2 and uses a deterministic encoder-decoder.

### B. Discovery of the Latent Directions

In this method, we use the latent representations produced by the trained DLVMs to discover latent directions representing ground truth factors. Given the *known* factors of variations,  $\mathcal{F}$ , for a dataset we use samples from the observed data corresponding to the  $i$ -th factor of variation,  $\mathcal{F}_i$ , to determine the direction in the latent space representing the latent variable for the factor,  $\mathcal{F}_i$ . To determine the direction corresponding to a factor  $\mathcal{F}_i$ ,  $L$  observed data are selected, where the factor  $\mathcal{F}_i$  is fixed to an *unique value* (chosen at random from  $\mathcal{V}_i$  containing the possible values for the  $i$ -th latent factor) for all  $L$  samples. The remaining factors,  $\mathcal{F}_{-i}$ , are assigned different values in each instance of the  $L$  examples, resulting in  $\mathcal{S}^i$  that is used to get the observed dataset,  $\mathcal{X}^i$ . The principal component analysis (PCA) is done on the encoding of  $L$  observed data in  $\mathcal{X}^i$ , produced by a trained DLVM, and the eigenvector with minimum variance,  $u_i$ , is chosen as the representative of the ground truth factor,  $\mathcal{F}_i$  (assigned a fixed value). This step is similar to the data generation technique of the FactorVAE metric [17] that associates a latent axis with a generative factor based on the variance along the latent axes.

Determination of the minimum variance eigenvector is repeated multiple times ( $N$ ) to capture the variation in  $u_i$ 's for different configurations of  $\mathcal{F}_i$  and  $\mathcal{F}_{-i}$ . The optimum unit vector,  $u^*$ , representing the factor,  $\mathcal{F}_i$ , is obtained by solving the optimization problem  $\max \sum_{i=1}^N (u^{*T} u_i)^2$ , where the  $u_i$  is an eigenvector. The solution to the optimization problem is the eigendecomposition of the mean outer product of the eigenvectors ( $u_i$ ), and  $u^*$  is the eigenvector with the maximum variance. The above steps are repeated to get the

---

### Algorithm 1 : Determine latent directions for generative factors in DLVMs, in a general setup

---

**Input:** Trained encoder  $E_\phi$ , Latent factors  $\mathcal{F}$ , Values for the latent factors  $\mathcal{V}$  ( $\mathcal{V}_i$  has the values for the factor,  $\mathcal{F}_i$ ),  $L$  samples used for the PCA analysis of a generative factor ( $\mathcal{F}_i$  is set to a fixed value chosen from  $\mathcal{V}_i$  and all others factors,  $\mathcal{F}_{-i}$ , i.e.,  $\mathcal{F} \setminus i$ , are allowed to vary), Number of PCA analysis ( $N$ ) to determine the direction of the generative factor,  $\mathcal{F}_i$ .

**Output:** Directions in the latent space,  $\mathcal{D}$ , for all the generative factors,  $\mathcal{F}$ .

```

1:  $\mathcal{D} \leftarrow \emptyset$ 
2: for Each ground truth factor  $k$  in  $\mathcal{F}$  do
3:   Sample  $N$  values for the ground truth factor  $k$ ,  $\mathcal{S}^k$ 
4:    $\mathcal{U} \leftarrow \emptyset$ 
5:   for each element  $\mathcal{S}_i^k$  in  $\mathcal{S}^k$  do
6:     Sample factors other than  $k$  ( $\mathcal{F}_{-k}$ )  $L$  times,  $\mathcal{S}^{-k}$ 
7:     {Concatenate  $\mathcal{S}_i^k$  to  $L$  samples in  $\mathcal{S}^{-k}$ }
8:      $\mathcal{S}^i \leftarrow \mathcal{S}^{-k} \cap \mathcal{S}_i^k$ 
9:     Get observed data,  $\mathcal{X}^i$ , corresponding to  $\mathcal{S}^i$ , where
10:     $\mathcal{X}^i = \{x_1^i, x_2^i \dots x_L^i\}$  and  $x_j^i \in \mathbb{R}^d$ 
11:     $\mathcal{Z}^i \leftarrow E_\phi(\mathcal{X}^i)$ , where  $\mathcal{Z}^i = \{z_1^i, z_2^i \dots z_L^i\}$  and  $z_j^i \in \mathbb{R}^l$ 
12:     $\{(\sigma_j, u_j)\}_{j=1}^l \leftarrow \text{PCA}(\mathcal{Z}^i)$ , where  $(\sigma_j, u_j)$  represents the eigenvector ( $u_j$ ) and the corresponding variance ( $\sigma_j$ ) estimated from the PCA
13:     $u_i$  is the eigenvector with the minimum variance
14:     $\mathcal{U} \leftarrow \mathcal{U} \cup u_i$ 
15:   end for
16:   {Comment: Estimate  $u^*$  representing the factor,  $\mathcal{F}_i$ }
17:    $\hat{\mathcal{U}} \leftarrow \mathbf{0}$ 
18:   for  $u_i$  in  $\mathcal{U}$  do
19:      $\hat{\mathcal{U}} \leftarrow \hat{\mathcal{U}} + u_i u_i^T$ 
20:   end for
21:    $\hat{\mathcal{U}} \leftarrow \frac{\hat{\mathcal{U}}}{N}$ 
22:    $\{(\sigma_j, u_j)\}_{j=1}^l \leftarrow \text{EIGEN}(\hat{\mathcal{U}})$ 
23:    $u^*$  is the eigenvector with the maximum variance
24:    $\mathcal{D} \leftarrow \mathcal{D} \cup u^*$ 
25: end for

```

---

latent directions ( $u^*$ ) for all the ground truth factors. The outline of the algorithm is presented in Algorithm 1.

### C. Evaluation

In this work, we follow the strategy proposed in the FactorVAE metric [17] and MIG metric [18] to devise techniques for disentanglement analysis using the latent directions,  $\mathcal{D}$ . However, using the estimated latent directions is not limited to these metrics. We name the proposed metrics as the *PCA FactorVAE metric* and the *PCA MIG metric*. Finding latent directions (representing latent variables),  $\mathcal{D}$ , corresponding to generative factors,  $\mathcal{F}$ , is a generalization of the majority vote classifier used in the FactorVAE metric. Thus, we can use the latent directions  $\mathcal{D}$  and  $u_i$  (estimated using the inner loop in algorithm 1) for different values of  $\mathcal{F}_i$  to predict the latent factor,  $\hat{\mathcal{F}}_i$ . In the PCA FactorVAE metric, we use a similarity

TABLE I

DISENTANGLEMENT SCORES OF COMPETING METHODS TRAINED WITH 10 DIFFERENT SEEDS FOR MULTIPLE DATASETS (HIGHER IS BETTER). THE BEST SCORE IS IN **BOLD**, AND THE SECOND BEST SCORE IS UNDERLINED. WE INDICATE THE IMPROVEMENT IN THE METRIC SCORES USING THE BLUE COLOR AND THE DROP WITH THE RED COLOR. WE OBSERVE THE MAXIMUM IMPROVEMENT IN THE METRIC SCORES OF THE AVAE FOR BOTH DATASETS.

Dataset	Method	FactorVAE $\uparrow$	PCA FactorVAE $\uparrow$	Diff $\uparrow$	MIG $\uparrow$	PCA MIG $\uparrow$	Diff $\uparrow$	MSE $\downarrow$
DSprites	VAE	64.78 $\pm$ 8.05	75.56 $\pm$ 7.21	<u>10.78</u>	0.06 $\pm$ 0.02	0.14 $\pm$ 0.04	<u>0.12</u>	3.68 $\pm$ 0.58
	$\beta$ -TCVAE	<b>75.55 <math>\pm</math> 3.52</b>	69.12 $\pm$ 15.08	<b>-6.43</b>	<b>0.20 <math>\pm</math> 0.06</b>	0.18 $\pm$ 0.13	<b>-0.02</b>	6.39 $\pm$ 2.05
	DIP-VAE-I	61.77 $\pm$ 8.96	70.68 $\pm$ 6.89	<u>8.91</u>	0.13 $\pm$ 0.07	0.12 $\pm$ 0.06	<b>-0.01</b>	3.61 $\pm$ 0.47
	DIP-VAE-II	60.70 $\pm$ 10.97	69.10 $\pm$ 2.15	<u>8.40</u>	0.08 $\pm$ 0.04	0.09 $\pm$ 0.02	<u>0.01</u>	3.46 $\pm$ 0.33
	RAE	64.21 $\pm$ 6.69	<b>82.03 <math>\pm</math> 2.58</b>	<u>17.82</u>	0.04 $\pm$ 0.01	0.16 $\pm$ 0.03	<u>0.12</u>	<b>2.51 <math>\pm</math> 0.20</b>
	WAE	47.87 $\pm$ 5.90	62.72 $\pm$ 6.85	<u>14.85</u>	0.02 $\pm$ 0.01	0.07 $\pm$ 0.02	<u>0.05</u>	3.69 $\pm$ 0.34
	AVAE	59.03 $\pm$ 2.62	<u>79.17 <math>\pm</math> 1.64</u>	<b>20.14</b>	0.02 $\pm$ 0.00	<b>0.20 <math>\pm</math> 0.02</b>	<b>0.18</b>	<u>2.98 <math>\pm</math> 0.28</u>
3D SHAPES	VAE	56.19 $\pm$ 8.85	55.00 $\pm$ 13.56	<b>-1.19</b>	0.13 $\pm$ 0.13	0.15 $\pm$ 0.15	<u>0.02</u>	10.47 $\pm$ 1.10
	$\beta$ -TCVAE	<b>75.51 <math>\pm</math> 12.84</b>	76.65 $\pm$ 15.78	<u>1.14</u>	<b>0.40 <math>\pm</math> 0.22</b>	0.46 $\pm$ 0.20	<u>0.06</u>	11.54 $\pm$ 1.86
	DIP-VAE-I	51.94 $\pm$ 1.91	48.94 $\pm$ 2.02	<b>-2.00</b>	0.06 $\pm$ 0.02	0.23 $\pm$ 0.03	<u>0.17</u>	<u>10.16 <math>\pm</math> 0.83</u>
	DIP-VAE-II	63.66 $\pm$ 11.26	66.04 $\pm$ 18.15	<u>2.38</u>	0.24 $\pm$ 0.18	0.27 $\pm$ 0.19	<u>0.03</u>	12.10 $\pm$ 2.64
	RAE	53.57 $\pm$ 13.14	70.85 $\pm$ 20.02	<u>17.28</u>	0.03 $\pm$ 0.02	0.33 $\pm$ 0.17	<u>0.30</u>	10.77 $\pm$ 1.25
	WAE	48.54 $\pm$ 3.04	52.34 $\pm$ 3.0	<u>3.80</u>	0.05 $\pm$ 0.03	0.20 $\pm$ 0.05	<u>0.15</u>	<b>9.80 <math>\pm</math> 1.95</b>
	AVAE	<u>72.90 <math>\pm</math> 7.30</u>	<b>91.93 <math>\pm</math> 3.27</b>	<b>19.03</b>	0.08 $\pm$ 0.02	<b>0.67 <math>\pm</math> 0.04</b>	<b>0.59</b>	10.29 $\pm$ 0.37

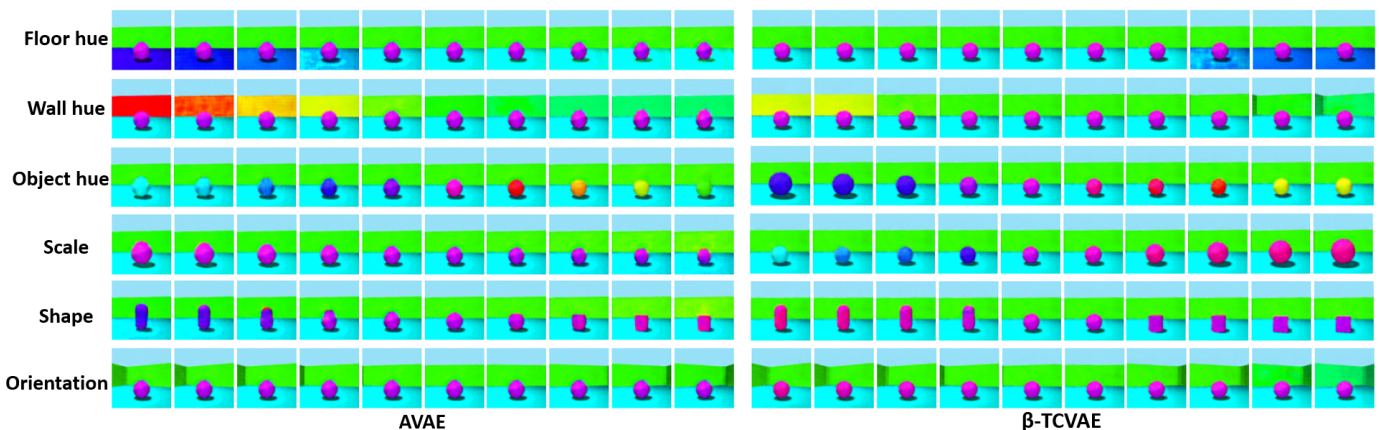


Fig. 1. Latent traversal of the 3D Shapes dataset [29] in the range  $[-\sigma, \sigma]$  for models trained using the AVAE and  $\beta$ -TCVAE. The latent factors are mentioned in the left column. All latent factors are represented by independent latent variables in the AVAE, with almost no overlap between latent variables except a slight variation in object color with shapes. For the  $\beta$ -TCVAE, we observe the entanglement of the multiple latent factors, such as the object color with the scale and the wall color with orientation. The visualization justifies the *improved* metrics scores of the AVAE in Table I using the proposed evaluation method.

measure between the  $u_i$  and the set of latent directions,  $\mathcal{D}$ , to predict the corresponding latent factor,  $\hat{\mathcal{F}}_i$ , and compare it to the true generative factor,  $\mathcal{F}_i$ , using cosine similarity measure (normalized correlation). The *prediction accuracy* of a model is the score for the PCA FactorVAE metric. In the PCA MIG metric, latent representations ( $\mathcal{Z} = E_\phi(\mathcal{X})$ ) are projected onto the latent directions,  $\mathcal{D}$ , and the MIG of the transformed representations ( $\mathcal{Z}' = \mathcal{Z}\mathcal{D}^T$ ) gives the score.

### III. EXPERIMENTS

#### A. Benchmark Methods & Datasets

Other than the regular VAE [1], we consider different variations of the VAE that modify the original formulation to match the aggregate posterior to the prior, such as the FactorVAE [17], and  $\beta$ -TCVAE [18], for comparison. We do not consider the FactorVAE as it uses a discriminator to optimize the objective function, which is challenging to train. For the same reason, we do not study the AAE [26] in this work. We study the DIP-VAE [19] that adds a regularizer to the VAE objective function to better match the aggregate posterior

to the prior. We also evaluate the RAE [30] and AVAE [5] in this work that uses a deterministic encoder, unlike the regular VAE. The AVAE is a new method based on the formulation of the VAE that addresses the posterior collapse and closely matches the aggregate posterior. Other than the variants of the VAE, we consider the WAE (with IMQ kernel) [27] that matches aggregate posterior in the latent space.

We use the DSprites [31] and 3D Shapes [29] datasets to evaluate DLVMs. The true generative factors of the observed data are known for both datasets. Annotated data is *required* for the quantification of the disentanglement in the latent representations of trained DLVMs.

#### B. Implementation Details

For a given dataset, we use the same latent dimension, encoder-decoder architecture (as used in [15]), and optimization strategies (such as the learning rate, learning rate scheduler, epochs, and batch size) for all the competing methods to ensure a fair comparison. We leverage the information of the known latent factors of the DSprites and 3D Shapes datasets

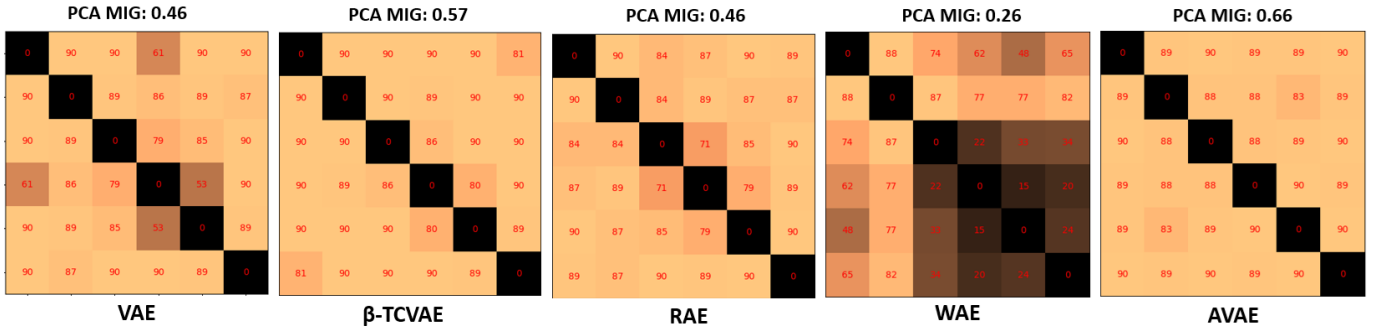


Fig. 2. Pairwise angle between the latent directions (six directions) estimated by Algorithm 1 for different DLVMs using the corresponding latent representations produced for the 3D Shapes dataset [29]. The latent directions should be orthogonal to each other for better disentanglement. Deviation from the orthogonality indicates entanglement of the ground truth generative factors that result in poor metric scores, as observed in the WAE. To interpret the estimated latent directions in the analysis of disentanglement, we report the corresponding PCA MIG metric score for each model.

TABLE II  
OPTIMIZATION SETTINGS FOR DIFFERENT METHODS.

Method	Parameters	DSprites	3D Shapes
$\beta$ -TCVAE	$\beta$ :	5	5
DIP-VAE-I	$(\lambda_{od}, \lambda_d)$ :	(10, 100)	(10, 100)
DIP-VAE-II	$(\lambda_{od}, \lambda_d)$ :	(10, 10)	(10, 10)
RAE	$\beta$ :	$1e-04$	$1e-04$
RAE	DEC-L2-REG:	$1e-06$	$1e-06$
WAE	RECONS-SCALAR:	0.05	0.05
WAE	$\beta$ :	10	10
AVAE	KDE samples ( $m$ ):	10K	10K

to set the latent size as  $l = 6$  for both datasets. For the AVAE, the number of the KDE samples used is 10K for both the DSprites and 3D Shapes datasets. We run all methods with 10 different seeds (producing different initialization) for the DSprites and 3D Shapes datasets. The objective function of several methods studied in this work has hyperparameters related to the regularization losses tuned for different datasets. Mostly, we have used hyperparameter settings suggested by the author or recommended in the literature [15], [18]. The hyperparameters of the methods are reported in Table II.

### C. Results

The performance of the competing methods under the proposed disentanglement metrics is reported in Table I for the DSprites [31] and 3D Shapes [29] datasets. A relatively poor reconstruction loss indicates stronger regularization of the latent representation, possibly leading to better disentanglement. Therefore, knowing the reconstruction loss of DLVMs on different datasets is informative. In an ideal scenario, we expect a higher disentanglement score with low reconstruction loss. In this experiment, we consider the FactorVAE metric [17] and the MIG metric [18] as the baseline to demonstrate the effectiveness of the proposed evaluation technique in Section II-C using the latent directions,  $\mathcal{D}$ , estimated by Algorithm 1.

Comparing the PCA MIG metric with the MIG metric [18] (baseline) demonstrates the impact of using the latent directions,  $\mathcal{D}$ , as latent variables relative to the latent axes. We observe an overall improvement in the performance of all the

methods studied in this work using the latent directions,  $\mathcal{D}$ , in the computation of the MIG scores, but the marginal drop on the DSprites dataset for the  $\beta$ -TCVAE and DIP-VAE-I. Likewise, except for the drop in the performance of the  $\beta$ -TCVAE on the DSprites dataset using the estimated latent directions, we observe a consistent improvement in the performance of all the methods. Methods matching the aggregate distribution benefit the most using the latent directions,  $\mathcal{D}$ .

Considering the metric scores reported in Table I, the AVAE produces the best score under both the metrics for the 3D Shapes dataset with a slightly higher MSE score (second best). This indicates that the AVAE achieves better disentanglement without compromising the quality of the reconstructed data, a desired property of a DLVM. The AVAE achieves significantly higher scores than the  $\beta$ -TCVAE for both the metrics and sets a new SOTA result for the 3D Shapes dataset. Higher metric scores for the AVAE are corroborated by the latent traversal along the latent directions,  $\mathcal{D}$ , shown in Figure 1. Overall, the performance of the AVAE is consistent under both metrics relative to the competing methods for both datasets.

The latent directions estimated by Algorithm 1 should be orthogonal for independent latent variables. The collapse of the latent directions or small angle between them indicates the entanglement of ground truth generative factors to latent variables that should affect the disentanglement scores. Figure 2 shows the angles between the estimated latent directions for different DLVMs trained on the 3D Shapes dataset. Every latent direction estimated for the AVAE is almost perpendicular to all other latent directions. This property explains the high metric scores of the AVAE in Table I on both datasets.

## IV. CONCLUSION

We present a statistical method to *evaluate disentanglement* in trained DLVMs using the estimated latent directions representing the generative factors of a dataset that *can be different from the cardinal latent axes*. We demonstrate improvements in metric scores for methods matching aggregate posterior distributions, such as the AVAE [5]. Therefore, we show limitations in the existing metrics that rely on cardinal latent axes

representing the generative factors. This work is supported by the National Institutes of Health grant R01ES032810.

## REFERENCES

- [1] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *International Conference on Learning Representations*, 2014.
- [2] D. J. Rezende, S. Mohamed, and D. Wierstra, “Stochastic backpropagation and approximate inference in deep generative models,” in *International Conference on Machine Learning*, 2014, pp. 1278–1286.
- [3] M. Rolinek, D. Zietlow, and G. Martius, “Variational autoencoders pursue pca directions (by accident),” in *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [4] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based generative modeling through stochastic differential equations,” in *International Conference on Learning Representations*, 2021.
- [5] S. Saha, S. Joshi, and R. Whitaker, “Matching aggregate posteriors in the variational autoencoder,” in *International Conference on Pattern Recognition*, 2024.
- [6] —, “Ard-vae: A statistical formulation to find the relevant latent dimensions of variational autoencoders,” 2025. [Online]. Available: <https://arxiv.org/abs/2501.10901>
- [7] Y. Bengio, A. Courville, and P. Vincent, “Representation learning: A review and new perspectives,” in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013.
- [8] M. Tschannen, O. Bachem, and M. Lucic, “Recent advances in autoencoder-based representation learning,” in *Workshop on Bayesian Deep Learning NeurIPS*, 2018.
- [9] S. Saha, O. Choi, and R. Whitaker, “Few-shot segmentation of microscopy images using gaussian process,” in *Medical Optical Imaging and Virtual Microscopy Image Analysis, MICCAI*, 2022.
- [10] S. Saha, W. Gazi, R. Mohammed, T. Rapstine, H. Powers, and R. Whitaker, “Multitask training as regularization strategy for seismic image segmentation,” *IEEE Geoscience and Remote Sensing Letters*, vol. 20, pp. 1–5, 2023.
- [11] X. Li, T. Mangin, S. Saha, E. Mohammed, Rehman Blanchard, D. Tang, H. Poppe, O. Choi, K. Kelly, and R. Whitaker, “Real-time idling vehicles detection using combined audio-visual deep learning,” in *In Emerging Cutting-Edge Developments in Intelligent Traffic and Transportation Systems, pages*, vol. 70, 2024, p. 142–158.
- [12] X. Li, R. Mohammed, T. Mangin, S. Saha, R. T. Whitaker, K. E. Kelly, and T. Tasdizen, “Joint audio-visual idling vehicle detection with streamlined input dependencies,” 2024. [Online]. Available: <https://arxiv.org/abs/2410.21170>
- [13] S. Ravi and H. Larochelle, “Optimization as a model for few-shot learning,” in *International Conference on Learning Representations*, 2017.
- [14] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International Conference on Machine Learning*, vol. 70, 2017, pp. 1126–1135.
- [15] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem, “Challenging common assumptions in the unsupervised learning of disentangled representations,” in *International Conference on Machine Learning*, 2019.
- [16] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, “ $\beta$ -vae: Learning basic visual concepts with a constrained variational framework,” in *International Conference on Learning Representations*, 2017.
- [17] H. Kim and A. Mnih, “Disentangling by factorising,” in *International Conference on Machine Learning*, 2018.
- [18] R. T. Q. Chen, X. Li, R. Grosse, and D. Duvenaud, “Isolating sources of disentanglement in vaes,” in *Conference on Neural Information Processing Systems*, 2019.
- [19] A. Kumar, P. Sattigeri, and A. Balakrishnan, “Variational inference of disentangled latent concepts from unlabeled observations,” in *International Conference on Learning Representations*, 2018.
- [20] C. Eastwood and C. K. I. Williams, “A framework for the quantitative evaluation of disentangled representations,” in *International Conference on Learning Representations*, 2018.
- [21] C. Eastwood, A. L. Nicolicioiu, J. V. Kügelgen, A. Kekić, F. Träuble, A. Dittadi, and B. Schölkopf, “DCI-ES: An extended disentanglement framework with connections to identifiability,” in *The Eleventh International Conference on Learning Representations*, 2023.
- [22] M. D. Hoffman and M. J. Johnson, “Elbo surgery: yet another way to carve up the variational evidence lower bound,” in *NIPS Workshop: Advances in Approximate Bayesian Inference*, 2016.
- [23] J. Lucasz, G. Tuckery, R. Grosse, and M. Norouzi, “Understanding posterior collapse in generative latent variable models,” in *International Conference on Learning Representations*, 2019.
- [24] J. R. Lucas, G. Tucker, R. B. Grosse, and M. Norouzi, “Don’t blame the elbo! a linear vae perspective on posterior collapse,” in *Conference on Neural Information Processing Systems*, 2019.
- [25] M. Rosca, B. Lakshminarayanan, and S. Mohamed, “Distribution matching in variational inference,” *arxiv*, 2018, preprint at <https://arxiv.org/abs/1802.06847>.
- [26] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, “Adversarial autoencoders,” in *International Conference on Learning Representations*, 2016.
- [27] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf, “Wasserstein auto-encoders,” in *International Conference on Learning Representations*, 2018.
- [28] S. Saha, S. Elhabian, and R. Whitaker, “Gens: generative encoding networks,” *Machine Learning*, vol. 111, p. 4003–4038, 2022.
- [29] C. Burgess and H. Kim, “3d shapes dataset,” <https://github.com/deepmind/3d-shapes/>, 2018.
- [30] P. Ghosh, M. S. M. Sajjadi, A. Vergari, M. Black, and B. Schölkopf, “From variational to deterministic autoencoders,” in *International Conference on Learning Representations*, 2020.
- [31] L. Matthey, I. Higgins, D. Hassabis, and A. Lerchner, “dsprites: Disentanglement testing sprites dataset,” <https://github.com/deepmind/dsprites-dataset/>, 2017.