



VAE 4: Challenging Assumptions behind Disentanglement

Shengyu Feng

Sep 28, 2021

Disentanglement review

Disentanglement is impossible without inductive bias

Experimental results

Future directions

Disentanglement review

Disentanglement is impossible without inductive bias

Experimental results

Future directions

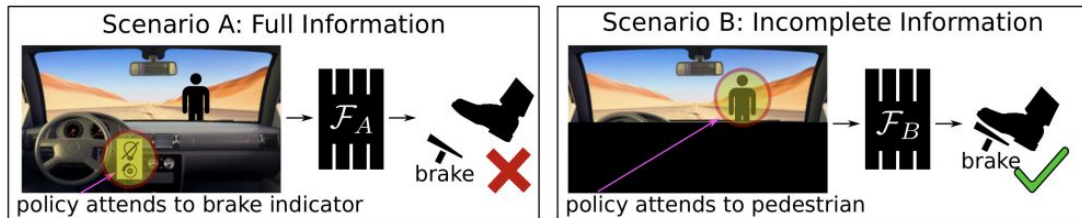
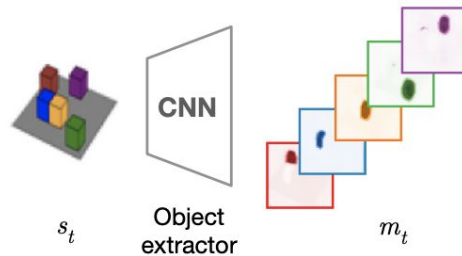
No formal definition, briefly, separating the distinct, informative factors of variations in the data

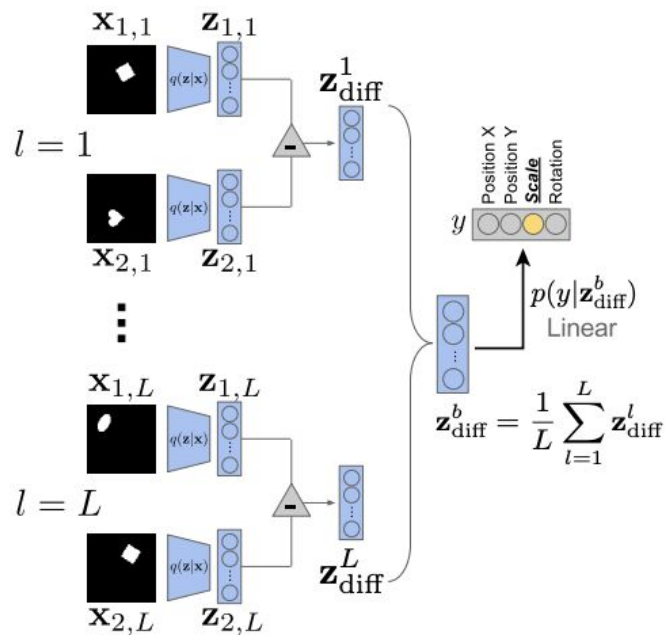
Entanglement: a random variable $X = [a+b, a-b, a+2b, b]$

Disentanglement: $Z = [a, b]$

Why is it useful?

- Customized generation
- Causality: $p(x)p(y|x), p(y)p(x|y)$
- Robot tasks





$$\mathcal{F}(\theta, \phi, \beta; \mathbf{x}, \mathbf{z}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Disentanglement review

Disentanglement is impossible without inductive bias

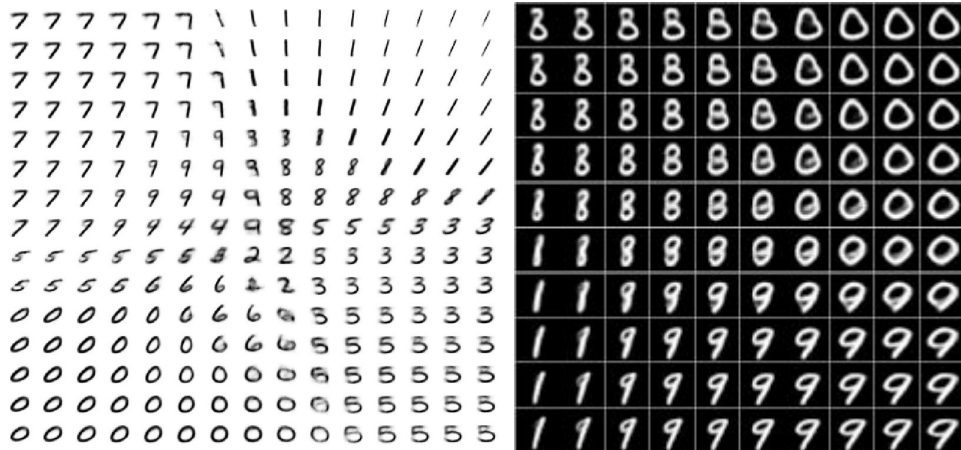
Experimental results

Future directions

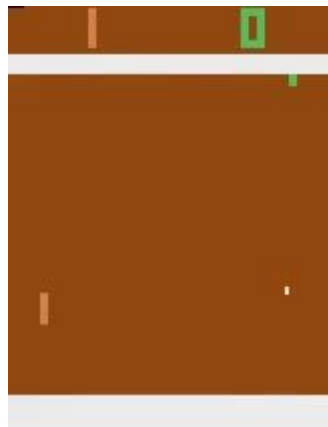
For any disentanglement representation, we can find infinite many equivalent representations. It's impossible to get the disentanglement representation without inductive bias.

Prior knowledge

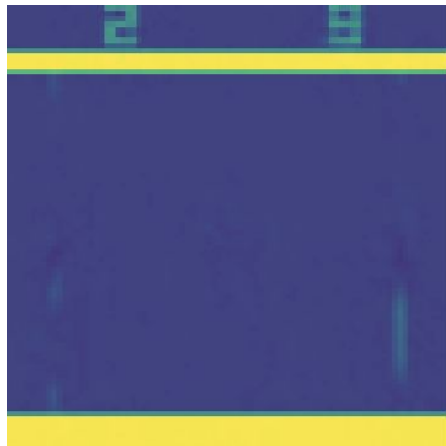
Theorem 1. For $d > 1$, let $\mathbf{z} \sim P$ denote any distribution which admits a density $p(\mathbf{z}) = \prod_{i=1}^d p(z_i)$. Then, there exists an infinite family of bijective functions $f : \text{supp}(\mathbf{z}) \rightarrow \text{supp}(\mathbf{z})$ such that $\frac{\partial f_i(\mathbf{u})}{\partial u_j} \neq 0$ almost everywhere for all i and j (i.e., \mathbf{z} and $f(\mathbf{z})$ are completely entangled) and $P(\mathbf{z} \leq \mathbf{u}) = P(f(\mathbf{z}) \leq \mathbf{u})$ for all $\mathbf{u} \in \text{supp}(\mathbf{z})$ (i.e., they have the same marginal distribution).



Data



Generated: $z = [0, 0, \dots, 0]$



Generated: $z = [0, 0, 0, 1, \dots, 0]$





Disentanglement review

Disentanglement is impossible without inductive bias

Experimental results

Future directions

Methods:

- β -VAE
- AnnealedVAE
- FactorVAE
- β -TCVAE
- DIP-VAE-I
- DIP-VAE-II

Metrics:

- β -VAE score
- FactorVAE score
- Mutual Information Gap (MIG)
- DCI Disentanglement
- Modularity
- SAP score

Datasets

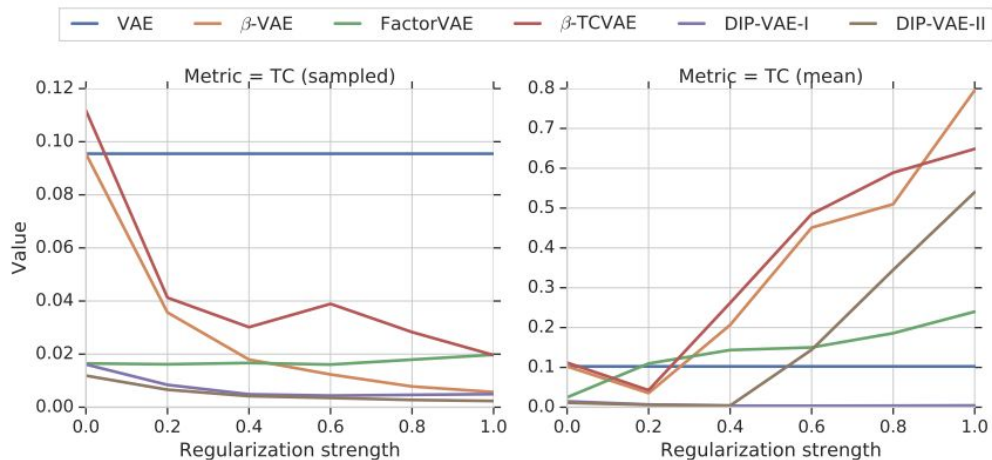
- dSprites
- Cars3D
- SmallNORB
- Shapes3D
- Color-dSprites
- Noisy-dSprites
- Scream-dSprites

Mean representation of the latent variables are correlated



Common practice in latent representation:

- Training: $z \sim N(\mu(z), \sigma(z))$
- Testig: $\mu(z)$



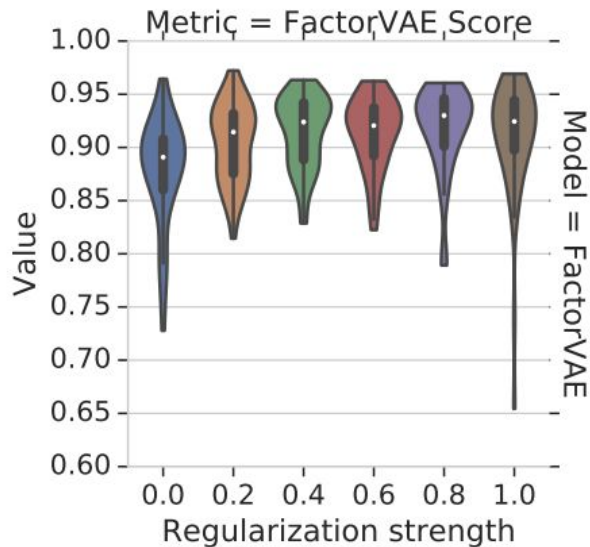
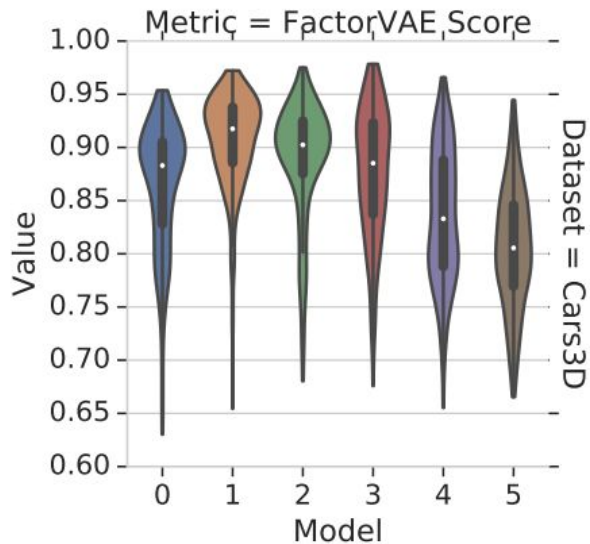
Most disentanglement metrics are correlated



Dataset = Noisy-dSprites

BetaVAE Score (A)	100	80	44	41	46	37
FactorVAE Score (B)	80	100	49	52	25	38
MIG (C)	44	49	100	76	6	42
Disentanglement (D)	41	52	76	100	-8	38
Modularity (E)	46	25	6	-8	100	13
SAP (F)	37	38	42	38	13	100
	(A)	(B)	(C)	(D)	(E)	(F)

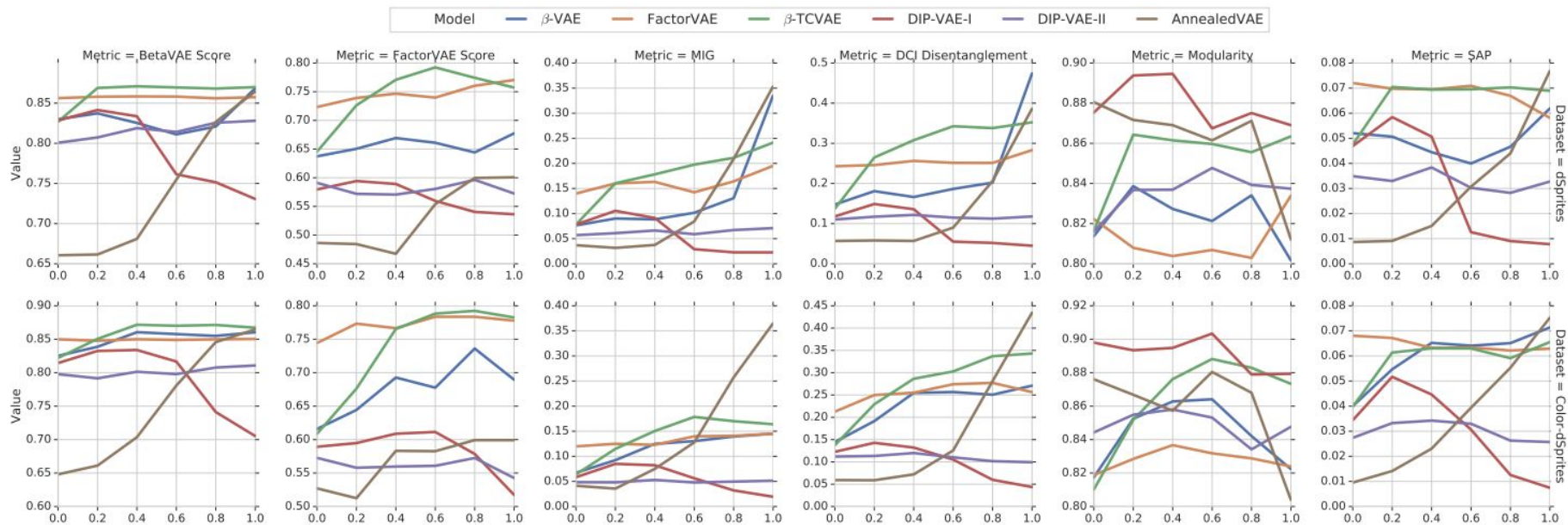
Hyperparameters and random seeds are important than models



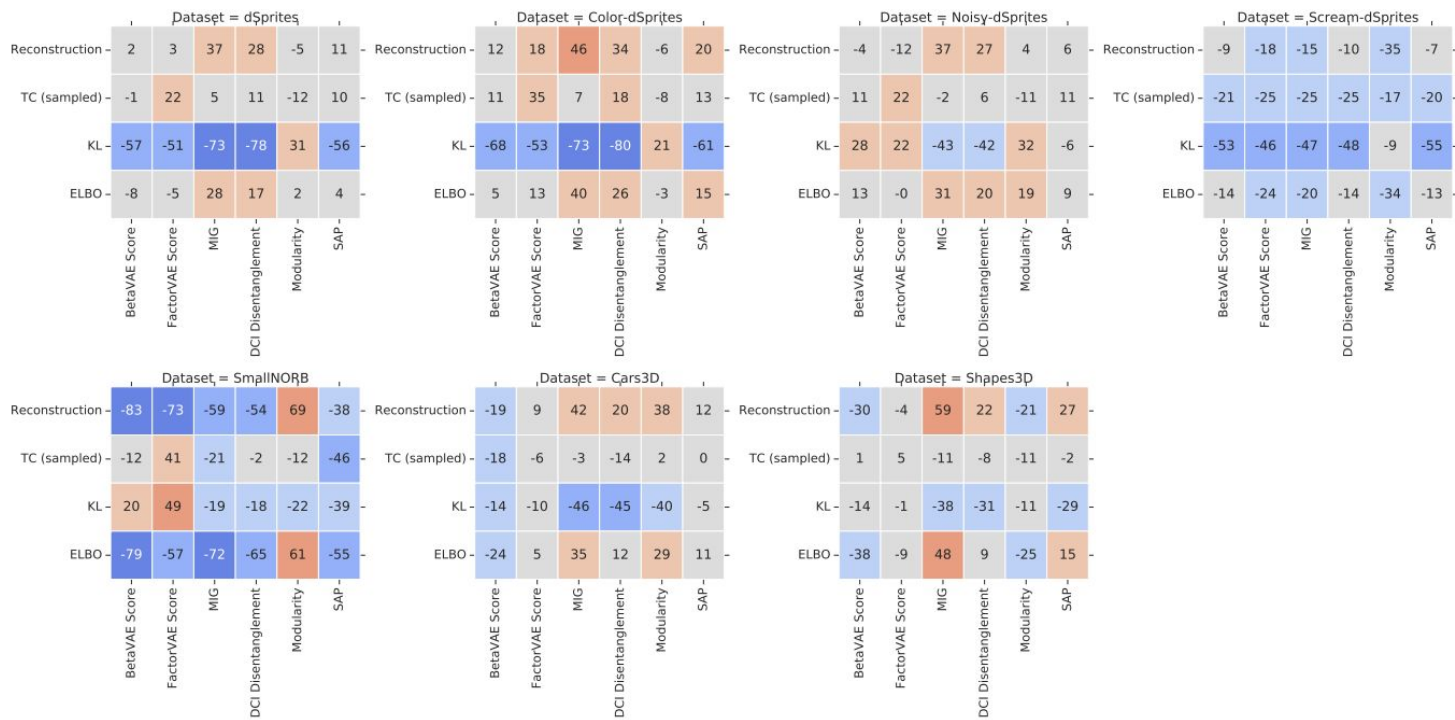
Recipes for hyperparameter selection



1. Strategy should not depend on the score, which needs labels and the full generative model
2. No hyperparameters and models work the best every time



Unsupervised loss vs. disentanglement scores



Bad transfer performance across datasets and metrics



The same dataset and metric: 80.7%

Different datasets and the same metric: 59.3%

Different datasets and metrics: 54.9%

Metric = DCI Disentanglement

dSprites (I)	100	95	65	65	34	64	46
Color-dSprites (II)	95	100	61	60	21	63	47
Noisy-dSprites (III)	65	61	100	68	17	64	59
Scream-dSprites (IV)	65	60	68	100	36	93	69
SmallNORB (V)	34	21	17	36	100	21	-9
Cars3D (VI)	64	63	64	93	21	100	85
Shapes3D (VII)	46	47	59	69	-9	85	100
	(I)	(II)	(III)	(IV)	(V)	(VI)	(VII)

Downstream tasks grounded on the disentangled representation

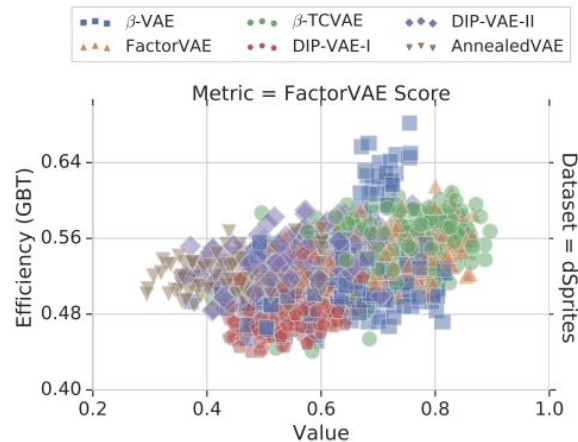


Basically, better disentanglement score leads to better performance in downstream tasks

But no clear evidence it leads to better sample complexity

Dataset = dSprites

	LR10	LR100	LR1000	LR10000	GBT10	GBT100	GBT1000	GBT10000	Efficiency (LR)	Efficiency (GBT)
BetaVAE Score	18	65	28	28	67	78	75	76	50	50
FactorVAE Score	13	49	13	12	58	73	71	71	43	46
MIG	18	63	20	-1	71	86	86	87	62	47
DCI Disentanglement	19	65	18	4	75	94	94	94	62	54
Modularity	-3	-9	15	18	-6	-17	-19	-13	-19	-14
SAP	12	64	20	12	71	77	74	75	56	49





Disentanglement review

Disentanglement is impossible without inductive bias

Experimental results

Future directions



Make inductive bias explicit, figure out how to select the hyperparameters without labels

Concrete the benefits of the disentanglement

Evaluate disentanglement methods over diverse datasets

- 2017 (ICLR): I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, A. Lerchner. [beta-VAE: Learning basic visual concepts with a constrained variational framework](#). ICLR, 2017.
- 2019 (ICML): F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Scholkopf, O. Bachem. [Challenging common assumptions in the unsupervised learning of disentangled representations](#). ICML, 2019.
- 2020 (ICLR): T. Kipf, E. Pol, M. Welling. [Contrastive Learning of Structured World Models](#). ICLR, 2020.
- 2019 (NeurIPS): P. Haan, D. Jayaraman, S. Levine. [Causal Confusion in Imitation Learning](#). NeurIPS, 2019.