# Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations

Francesco Locatello <sup>12</sup> Stefan Bauer <sup>2</sup> Mario Lucic <sup>3</sup> Gunnar Rätsch <sup>1</sup> Sylvain Gelly <sup>3</sup> Bernhard Schölkopf <sup>2</sup> Olivier Bachem <sup>3</sup>

- 1. ETH Zurich, Department for Computer Science
- 2. Max Planck Institute for Intelligent Systems
- 3. Google Research, Brain Team

**ICML 2019** 

Hanqing Chao

# Representation Learning Problem Definition

#### Step 1

- A multivariate latent random variable z is sampled from a distribution P(z).
- Intuitively, **z** corresponds to **semantically meaningful factors** of variation of the observations (e.g., content + position of objects in an image).

### Step 2

• The observation  $\mathbf{x}$  is sampled from the conditional distribution  $P(\mathbf{x}|\mathbf{z})$ .

#### Goal

• Find useful transformations  $r(\mathbf{x})$  of  $\mathbf{x}$  that "make it easier to extract useful information when building classifiers or other predictors" (Bengio et al., 2013).



# Definition of Disentanglement informal

A change in a single underlying factor of variation  $z_i$  should lead to a change in a single factor in the learned representation r(x).

# Frequency VS Bayesian

**Maximum Likelihood Estimation** 

$$p(x^{(1)},\cdots,x^{(m)}| heta)$$

# Frequency VS Bayesian



#### **Maximum Posterior Estimation**

$$p( heta|x^{(1)},\cdots,x^{(m)}) = rac{p(x^{(1)},\cdots,x^{(m)}| heta)p( heta)}{p(x^{(1)},\cdots,x^{(m)})}$$

#### On Test

$$p(x^{(m+1)}|x^{(1)},\cdots,x^{(m)}) = \int p(x^{(m+1)}| heta)p( heta|x^{(1)},\cdots,x^{(m)})d heta$$

## VAE

- Assume a specific prior P(z) on the latent space.
- Parameterize the conditional probability  $P(\mathbf{x}|\mathbf{z})$  using a DNN. (Decoder)
- The distribution  $P(\mathbf{z}|\mathbf{x})$  is approximated using a variational distribution  $Q(\mathbf{z}|\mathbf{x})$  ( $Q(\mathbf{z}|\mathbf{x}_i)$ ). (Encoder)

The representation for  $r(\mathbf{x})$  is usually taken to be the mean of the approximate posterior distribution  $Q(\mathbf{z}|\mathbf{x})$ .



# Disentanglement

Try to enforce a factorized aggregated

posterior 
$$\int_{\mathbf{x}} Q(\mathbf{z}|\mathbf{x})P(\mathbf{x}) d\mathbf{x}$$
,

**Theorem 1.** For d > 1, let  $\mathbf{z} \sim P$  denote any distribution which admits a density  $p(\mathbf{z}) = \prod_{i=1}^d p(\mathbf{z}_i)$ . Then, there exists an infinite family of bijective functions  $f: \operatorname{supp}(\mathbf{z}) \to$  $\operatorname{supp}(\mathbf{z})$  such that  $\frac{\partial f_i(\mathbf{u})}{\partial u_i} \neq 0$  almost everywhere for all i and j (i.e., **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).

- Assume we have p(z) and some P(x|z) defining a generative model.
- Consider any unsupervised disentanglement method and assume that it finds a representation r(x) that is perfectly disentangled with respect to z in the generative model.
- Theorem 1 implies that there is an equivalent generative model with the latent variable  $z^{\hat{}}=f(z)$  where  $z^{\hat{}}$  is completely entangled with respect to z and thus also r(x)
- Furthermore, since f is deterministic and  $p(z) = p(z^{\hat{}})$  almost everywhere, both generative models have the same marginal distribution of the observations x by construction, i.e.

$$P(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int p(\mathbf{x}|\hat{\mathbf{z}})p(\hat{\mathbf{z}})d\hat{\mathbf{z}}.$$

The (unsupervised) disentanglement method only has access to observations x, it hence cannot distinguish between the two equivalent generative models and thus has to be entangled to at least one of them.



- Assume we have p(z) and some P(x|z) defining a generative model.
- Consider any unsupervised disentanglement method and assume that it finds a representation r(x) that is perfectly disentangled with respect to z in the generative model.
- Theorem 1 implies that there is an equivalent generative model with the latent variable  $z^{\hat{}} = f(z)$  where  $z^{\hat{}}$  is completely entangled with respect to z and thus also r(x)
- Furthermore, since f is deterministic and  $p(z) = p(z^{\hat{}})$  almost everywhere, both generative models have the same marginal distribution of the observations x by construction, i.e.

$$P(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int p(\mathbf{x}|\hat{\mathbf{z}})p(\hat{\mathbf{z}})d\hat{\mathbf{z}}.$$

The (unsupervised) disentanglement method only has access to observations x, it hence cannot distinguish between the two equivalent generative models and thus has to be entangled to at least one of them.



- Assume we have p(z) and some P(x|z) defining a generative model.
- Consider any unsupervised disentanglement method and assume that it finds a representation r(x) that is perfectly disentangled with respect to z in the generative model.
- Theorem 1 implies that there is an equivalent generative model with the latent variable  $z^{\hat{}} = f(z)$  where  $z^{\hat{}}$  is completely entangled with respect to z and thus also r(x)
- Furthermore, since f is deterministic and  $p(z) = p(z^{\hat{}})$  almost everywhere, both generative models have the same marginal distribution of the observations x by construction, i.e.

$$P(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int p(\mathbf{x}|\hat{\mathbf{z}})p(\hat{\mathbf{z}})d\hat{\mathbf{z}}.$$

The (unsupervised) disentanglement method only has access to observations x, it hence cannot distinguish between the two equivalent generative models and thus has to be entangled to at least one of them.



- Assume we have p(z) and some P(x|z) defining a generative model.
- Consider any unsupervised disentanglement method and assume that it finds a representation r(x) that is perfectly disentangled with respect to z in the generative model.
- Theorem 1 implies that there is an equivalent generative model with the latent variable  $z^{\hat{}} = f(z)$  where  $z^{\hat{}}$  is completely entangled with respect to z and thus also r(x)
- Furthermore, since f is deterministic and  $p(z) = p(z^{\hat{}})$  almost everywhere, both generative models have the same marginal distribution of the observations x by construction, i.e.

$$P(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int p(\mathbf{x}|\hat{\mathbf{z}})p(\hat{\mathbf{z}})d\hat{\mathbf{z}}.$$

The (unsupervised) disentanglement method only has access to observations x, it hence cannot distinguish between the two equivalent generative models and thus has to be entangled to at least one of them.



We need **inductive biases** on **both** datasets and models



# Experiments

- 6 models
- 6 disentangle measurements
- 50 random initial
- 7 datasets
- 12,000 models
- https://github.com/google-research/disentanglement\_lib



# Experiments



dSprites is a dataset of 2D shapes procedurally generated from 6 ground truth independent latent factors. These factors are *color*, *shape*, *scale*, *rotation*, *x* and *y* positions of a sprite.

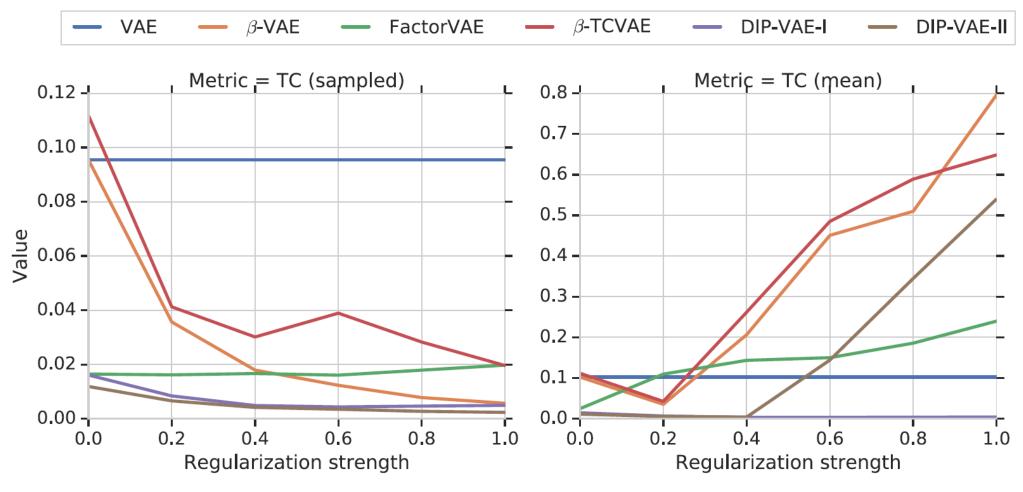
All possible combinations of these latents are present exactly once, generating N = 737280 total images.

#### Latent factor values

- Color: white
- Shape: square, ellipse, heart
- Scale: 6 values linearly spaced in [0.5, 1]
- Orientation: 40 values in [0, 2 pi]
- Position X: 32 values in [0, 1]
- Position Y: 32 values in [0, 1]

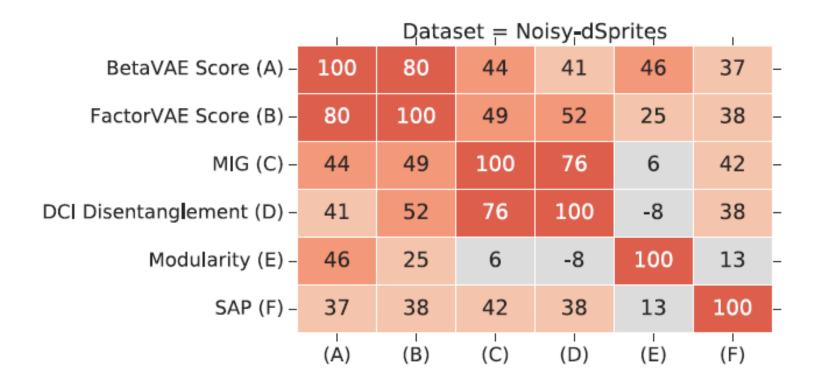


# Can current methods enforce a uncorrelated aggregated posterior and representation?



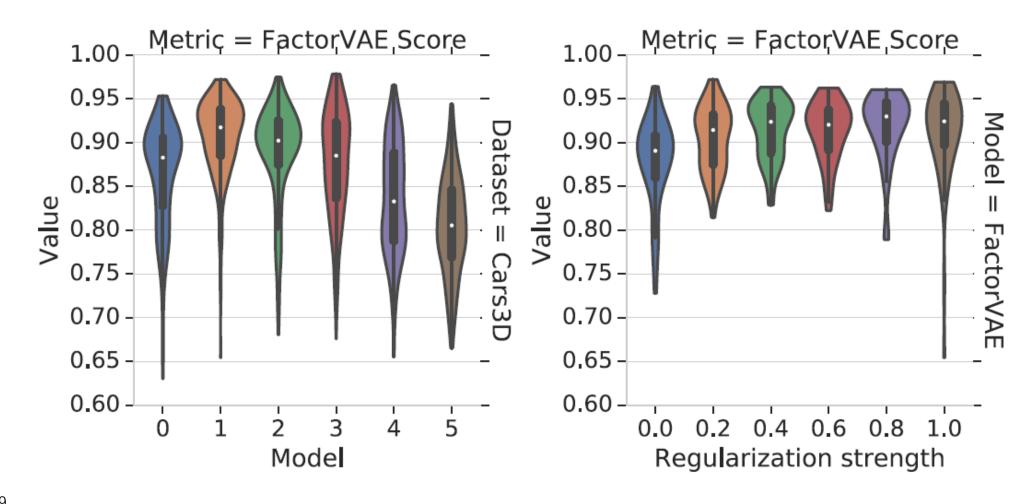


# How much do the disentanglement metrics agree?

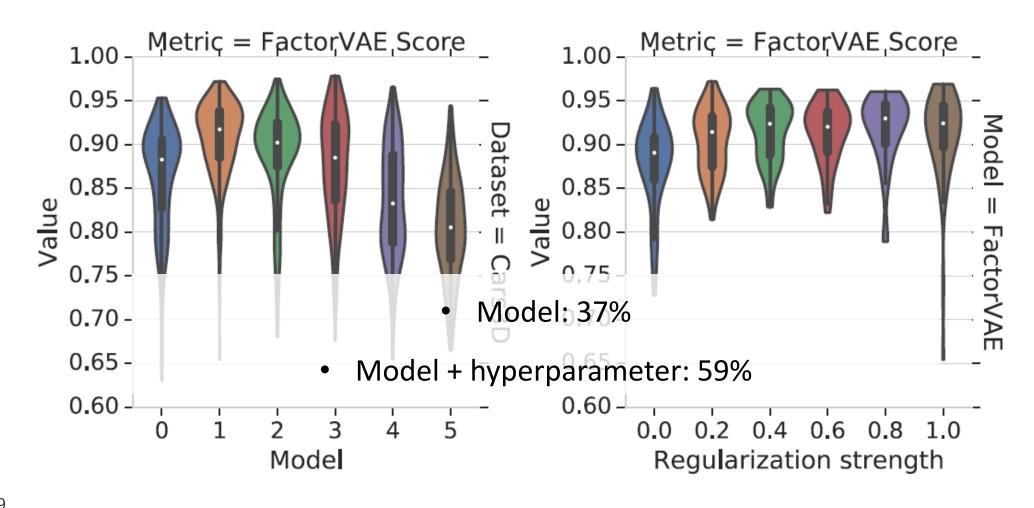




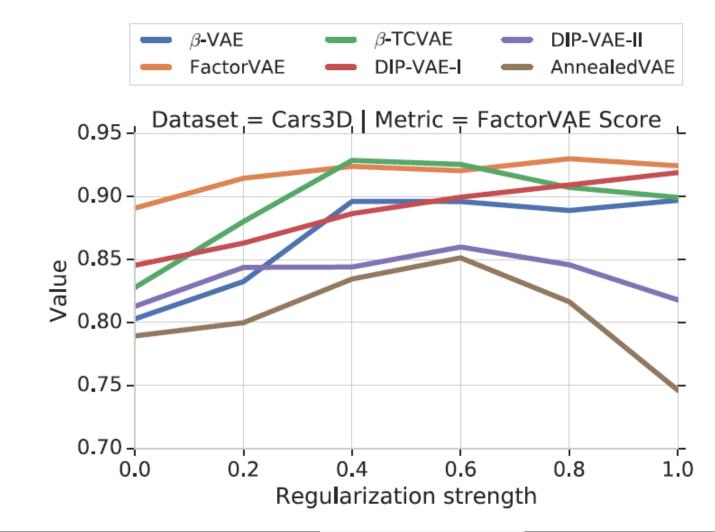
# How important are different models and hyperparameters for disentanglement?



# How important are different models and hyperparameters for disentanglement?



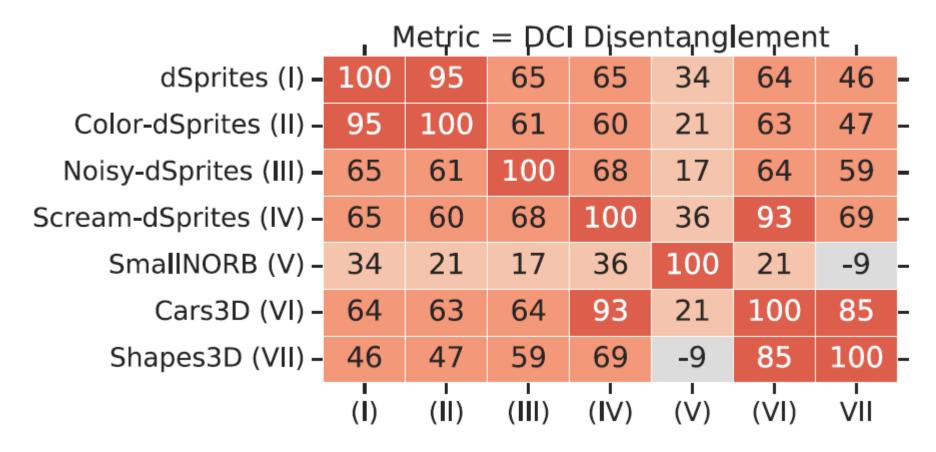
# Are there reliable recipes for model selection? General recipes for hyperparameter selection.



# Are there reliable recipes for model selection? *Model selection based on unsupervised scores.*



# Are there reliable recipes for model selection? Hyperparameter selection based on transfer.



# Are there reliable recipes for model selection?

Unsupervised model selection remains an unsolved problem.

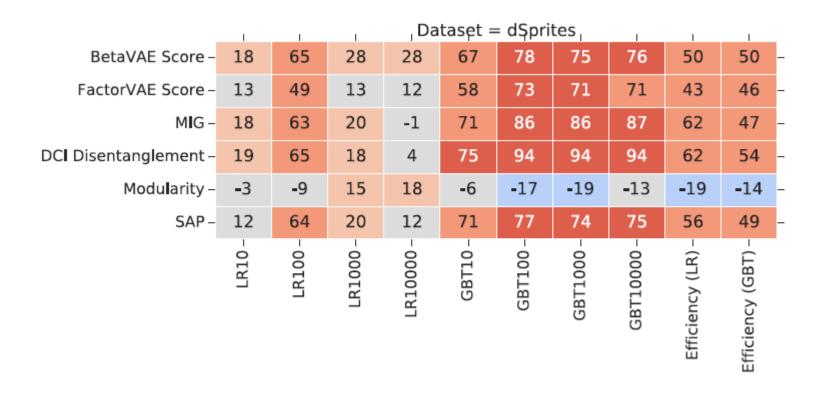


# Are these disentangled representations useful for downstream tasks in terms of the sample complexity of learning?

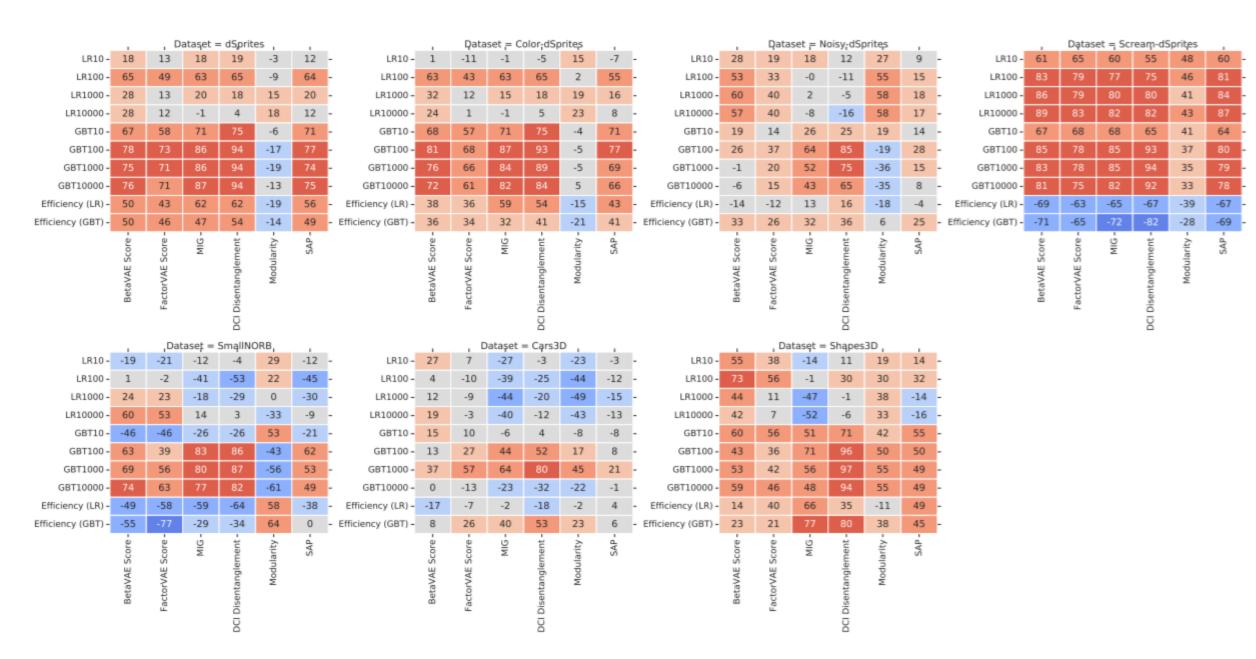
#### Goal

To **recover the true factors** of variations from the learned representation using either multi-class logistic regression (LR) or gradient boosted trees (GBT).

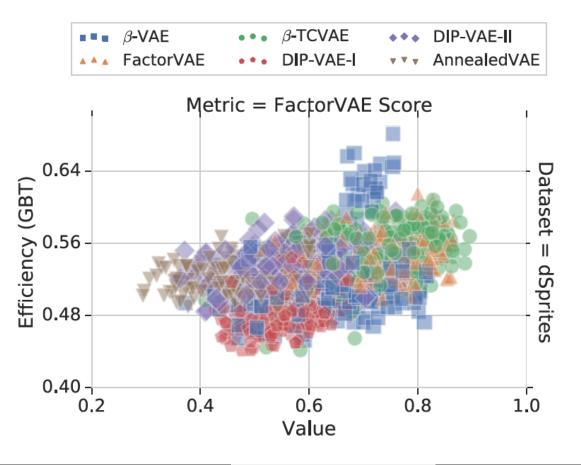
# Are these disentangled representations useful for downstream tasks in terms of the sample complexity of learning?

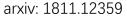






# Are these disentangled representations useful for downstream tasks in terms of the sample complexity of learning?





## Conclusion

- Inductive biases and implicit and explicit supervision. Deviates from the static, purely unsupervised setting.
- Concrete practical benefits of disentangled representations.
- Experimental setup and diversity of data sets. it is easy to draw spurious conclusions from experimental results if one only considers a subset of methods, metrics and data sets.

## Conclusion

- Inductive biases and implicit and explicit supervision. Deviates from the static, purely unsupervised setting.
- Concrete practical benefits of disentangled representations.
- Experimental setup and diversity of data sets. it is easy to draw spurious conclusions from experimental results if one only considers a subset of methods, metrics and data sets.

## Conclusion

- Inductive biases and implicit and explicit supervision. Deviates from the static, purely unsupervised setting.
- Concrete practical benefits of disentangled representations.
- Experimental setup and diversity of data sets. it is easy to draw spurious conclusions from experimental results if one only considers a subset of methods, metrics and data sets.

# Thanks !