Université m de Montréal

Sébastien Lachapelle¹, Pau Rodríguez López², Yash Sharma³, Katie Everett⁴, Rémi Le Priol¹, Alexandre Lacoste², Simon Lacoste-Julien^{1,5} ¹Mila, Université de Montréal ²ServiceNow Research ³Tübingen Al Center, University of Tübingen ⁴Google Reasearch ⁵Canada CIFAR Al Chair

CONTRIBUTIONS

- A new approach to achieve nonlinear ICA/disentanglement.
- Novel identifiability result based on the assumption that the causal graph connecting the latent factors is typically sparse [1, 2].
 - Objects usually interact sparsely with each others.
 - Actions typically affect only a few factors of variations.
- An estimation procedure which relies on the variational autoencoder (VAE) and a masking mechanism allowing to regularize for sparsity.
- As a special case of our framework, showing how unknown-target interventions on the latent factors can be leveraged to disentangle them, which is closely related to the **sparse mechanism shift hypothesis** introduced in [10].
- Illustrations of our theory with synthetic data.

- BACKGROUND

1.1 Motivating example

(Latent factors)

(T) (Tree position)

Z = |(R)| (Robot position)

(B) (Ball position)

f (Nonlinear decoder)



High-dim. observation

 $X = \mathbf{f}(Z)$

- $X \in \mathbb{R}^{d_x}$: High-dim. observation
- $Z \in \mathbb{R}^{d_z}$: Low-dim. latent representation
- $\mathbf{f}: \mathbb{R}^{d_z} \to \mathcal{X} \subset \mathbb{R}^{d_x}$: Invertible decoder
- $d_z \leq d_x$

1.2 Disentanglement

Definition (Disentanglement) Given a ground-truth model f, a learned representation f is *disentangled* if both are **permutation-equivalent**.

Definition (Permutation-equivalence)^{*a*} [7]

Two models f and \hat{f} (representing the same data manifold \mathcal{X}) are permutation*equivalent*, denoted by $\mathbf{f} \sim_p \mathbf{f}$, when

$$\mathbf{f}^{-1}(x) = PD\tilde{\mathbf{f}}^{-1}(x) + b \ \forall x \in \mathcal{X}$$

where P is a permutation matrix and D an invertible diagonal matrix. ^aSimplified definition, see paper for details.

1.3 Identifiability

• Unsupervised disentanglement is possible when f is identifiable up to \sim_p

 \sim_p -identifiability: $p(x) = \hat{p}(x) \ \forall x \implies \mathbf{f} \sim_p \hat{\mathbf{f}}$,

where p(x) and $\hat{p}(x)$ are the densities of X when the decoder is f and \hat{f} , respectively.

- **Problem:** With the standard assumption of factor independence, the mixing function f is **not identifiable** from the distribution of X (without further assumptions) [3].
- **Recent development:** Identifiability of the latent factors is possible even with nonlinear mixing, as long as the latent variables are **conditionally independent given** an observed auxiliary variable [4, 7, 6].
- $A \in \mathbb{R}^{d_a}$: Observed auxiliary variable, e.g. an environment index or an action.

Disentanglement via Mechanism Sparsity Regularization: A New Principle for Nonlinear ICA



