Learning Causal Representations Apprentissage de représentations causales



Simon Lacoste-Julien



DIRO, Université de Montréal & Mila



Montreal, April 20th 2023



Credits to:





Sébastien Lachapelle

Why causality in machine learning?

- ML researchers are becoming more concerned with **robustness** of their model to **out-of-distribution generalization**
 - One way: use causal models
 - Handle interventions



• E.g. in reinforcement learning, policy applications, etc. -> need to model the effect of interventions

Learning causal representations

 But what if you don't even know the identity of causal variables (high level semantic variables)?
 -> need to learn them!

Towards Causal Representation Learning

Bernhard Schölkopf[†], Francesco Locatello[†], Stefan Bauer^{*}, Nan Rosemary Ke^{*}, Nal Kalchbrenner Anirudh Goyal, Yoshua Bengio

Proceedings of the IEEE 2021

• Motivating example:



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

Sébastien Lachapelle¹, Pau Rodríguez López², Yash Sharma³, Katie Everett⁴, Rémi Le Priol¹, Alexandre Lacoste² and Simon Lacoste-Julien¹







Google Research

CLeaR 2022



Motivating Example

Disentanglement is the problem of recovering the latent factors without supervision, i.e. from p(x)

Our **identifiability theory** shows how to use this sparsity to disentangle the latent factors via **sparsity regularization**



The general problem of identifiability for generative models

Consider the following simple generative model:

$$Z \sim \mathbb{P}_{Z}, X := \mathbf{f}(Z) \implies \mathbb{P}_{X}$$
Consider this other model:

$$II \quad \text{Both models represent the same distribution over X...}$$

$$\widehat{I} := UZ, \hat{X} := \mathbf{f}(U^{-1}, \hat{Z}) \implies \mathbb{P}_{\hat{X}}$$
... but their representations can be drastically different $\widehat{\mathbf{f}}$

What is disentanglement? (Ground-truth) $\mathcal{Z} = \mathbb{R}^{d_z}$ $\hat{\mathcal{Z}} = \mathbb{R}^{d_{z \, \text{(Learned)}}}$ Model isnatsdiagentangled! $\mathbf{f}^{-1} \circ \hat{\mathbf{f}}$ (Learned decoder) (Ground-truth decoder) \mathbf{f} F (Data manifold) $\mathcal{X} \subset \mathbb{R}^{d_x}$ (Observation space)

Bird's-eye view of our approach

- Must learn
 - \circ a decoder **f**
 - a latent transition model $p(z^t \mid z^{< t}, a^{< t})$ -> parameter λ
 - \circ Causal graph over the latents/actions G
- Our theory shows how regularizing *G* to be sparse will help us identify the latent variables, i.e. to have a disentangled representation.



Learnable parameters: $\theta := (\mathbf{f}, \lambda, G)$

Useful to think of

- θ = parameter of the ground-truth model
- $\hat{\theta}$ or $\tilde{\theta}$ = parameter of the learned model

Our Theorem: Disentanglement via Mechanism Sparsity

(informally)

- If θ and $\hat{\theta}$ model the same distribution on X
- If the ground truth transition is "sufficiently complex"
- If we ensure our predicted graph is **as sparse** as the ground truth graph (this can be done via regularization)
- If the ground truth graph is **sufficiently sparse** (precise in paper)

Then θ and $\hat{\theta}$ are permutation-equivalent *i.e. the model* $\hat{\theta}$ *is disentangled.*

-> Generalized to no constraints on ground truth graph to get **block-permutation equivalence** in paper *"Partial Disentanglement via Mechanism Sparsity"*, [Lachapelle & Lacoste-Julien, arXiv 2022]

10

Our proposed method

- Model transition model and decoder via neural networks
- Use **binary masks** to encode transition graphs
- Use VAE approach to do approximate maximum likelihood on data
- Use **Gumbel-softmax trick** to learn discrete masks estimate gradients via reparameterization trick [Jang et al., 2017, Maddison et al., 2017]
 - Can also easily do I0-regularization (sparse graphs) with it

Synthetic Toy Experiments

 \bigcirc

 ∇



- R^2 = Measure of linear equivalence (higher is better)
 - mcc = Measure of permutation equivalence/disentanglement (higher is better)
 - shd = Measure graph recovery (lower is better)

Ongoing work: experiments on Atari games (e.g. pong)

Other work: causal graphical model learning when you know the causal variables

Differentiable Causal Discovery from Interventional Data (DCDI) NeurIPS 2020



Philippe Brouillard*



Sébastien Lachapelle*



Alexandre Lacoste



Simon Lacoste-Julien



Alexandre Drouin

Taxonomy of score-based algorithms (non-exhaustive)



Conclusion & outlook

- Early steps to learn causal variables from high dimensional data
 - By exploiting sparsity of interactions
- Still need to experiment with more realistic data
- Follow-up work showing the benefit of approach for **multi-task learning** with sparse parameters with sufficient variability across tasks
- Still space for algorithmic improvements, scaling and more theory (finite sample analysis,etc)
- Potential for model-based reinforcement learning

Thank you! – Merci!