

SAM

Structural Agnostic Model, causal discovery and penalized adversarial learning

D.Kalainathan, O.Goudet, D.Lopez-Paz, P.Caillou, I.Guyon, M.Sebag

TAU, CNRS, INRIA,
Université Paris Sud,
Université Paris Saclay, France
Facebook AI Research

Outline of the talk

Introduction

State of the art

Learning functional Causal Models with generative neural networks

SAM: Redefining causal graph recovery as a single optimization problem

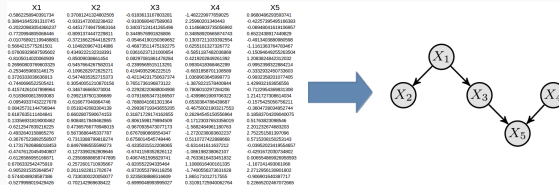
Experiments

Towards improving SAM

Introduction

Causal Discovery

- Goal: build a graph which models how the data could have been generated



- Gold Standard: perform randomized experiments
- Our setting: having only observational data: infer causal relationships in the dataset.

Correlation does not imply causation

- Task: predict a target variable Y given (X_1, X_2)
- Generative process underlying (X_1, X_2, Y) :

$$X_1, E_{X_1}, E_{X_2} \sim \text{Uniform}(0, 1), X_1 \perp\!\!\!\perp E_{X_1}, Y \perp\!\!\!\perp E_{X_2}$$

$$Y \leftarrow 0.5X_1 + E_{X_1},$$

$$X_2 \leftarrow Y + E_{X_2},$$



- Least-squares solution: $\hat{Y} = 0.25X_1 + 0.5X_2$
- X_2 is a better predictor for Y than X_1
- However X_2 does not cause Y

State of the art

Causality - Key idea 1 : Identifying v-structures

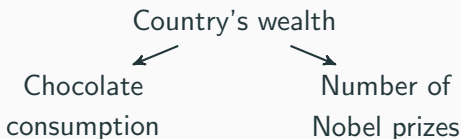
Exploit conditional independence to identify causal relations
[Spirtes et al., 2000], [Tsamardinos et al., 2006]:



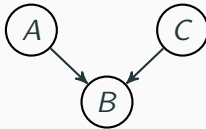
3 Markov Equivalent Classes: $A \perp\!\!\!\perp C | B$



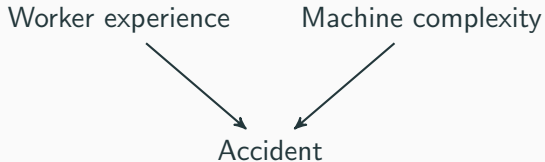
Example



V-Structure: $A \not\perp C | B$

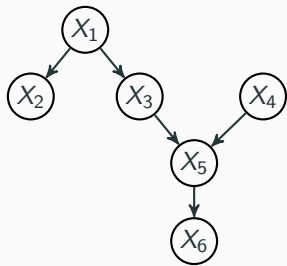


Example

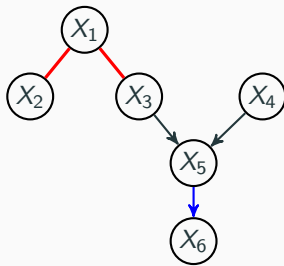


Constraint-based methods

Constraint-based methods, through V-Structures and constraint propagation, output a **CPDAG** (Completed Partially Directed Acyclic Graph). PC algorithm [Spirtes et al., 2000].



(a) The exact DAG of \mathcal{G} .



(b) The CPDAG of \mathcal{G} .

Limitations

- Explosive number of conditional independence tests to perform
- Cannot identify all cause effect relations

$$X_1, E_{X_1}, E_{X_2} \sim \text{Uniform}(0, 1), X_1 \perp\!\!\!\perp E_{X_1}, Y \perp\!\!\!\perp E_{X_2}$$

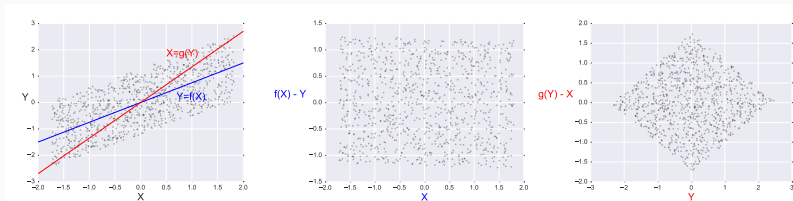
$$Y \leftarrow 0.5X_1 + E_{X_1},$$

$$X_2 \leftarrow Y + E_{X_2},$$



- Here $X_1 \perp\!\!\!\perp X_2 | Y$. No V-structure

Key idea 2: exploit asymmetry between cause and effect



- Causal additive noise model (ANM) [Hoyer et al., 2009]:
 $Y = f(X) + E$, with $X \perp\!\!\!\perp E$
- Perform a regression and check independence of the residual and the cause

- With the previous example and a linear causal additive noise model:

$$X_1, E_{X_1}, E_{X_2} \sim \text{Uniform}(0, 1), X_1 \perp\!\!\!\perp E_{X_1}, Y \perp\!\!\!\perp E_{X_2}$$

$$Y \leftarrow 0.5X_1 + E_{X_1},$$

$$X_2 \leftarrow Y + E_{X_2},$$

- Regression Y on X_1 : $X_1 \rightarrow Y$
- Regression X_2 on Y : $Y \rightarrow X_2$

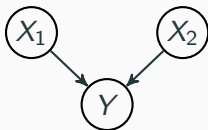


Limitation

- Generality problem: sometimes the causal additive noise model does not hold in any direction
e.g.: does not work for $Y = X \times E$
- Do not take into account independence relations. Consider:

$$X_1, X_2, E_{X_1} \sim \text{Gaussian}(0, 1), X_1 \perp\!\!\!\perp E_{X_1}, X_2 \perp\!\!\!\perp E_{X_1}$$

$$Y \leftarrow 0.5X_1 + X_2 + E_{X_1}$$



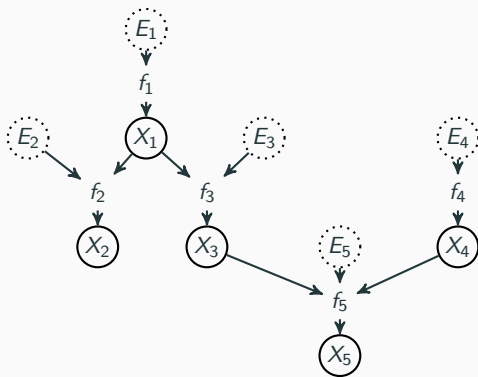
- (X_1, Y) and (X_2, Y) are perfect symmetric pairwise distribution (after rescaling)
- However $X_1 \not\perp\!\!\!\perp X_2 | Y$. A V-structure may be identified

Learning functional Causal Models with generative neural networks

Functional Causal Models (FCMs)

$$X_i = f_i(X_{\text{Pa}(i;\mathcal{G})}, E_i), \forall i \in [1, d]$$

$X_{\text{Pa}(i;\mathcal{G})}$ the set of parents of X_i in \mathcal{G} , E_i a random independent noise variable, f_i a deterministic function



$$\begin{cases} X_1 = f_1(E_1) \\ X_2 = f_2(X_1, E_2) \\ X_3 = f_3(X_1, E_3) \\ X_4 = f_4(E_4) \\ X_5 = f_5(X_3, X_4, E_5) \end{cases}$$

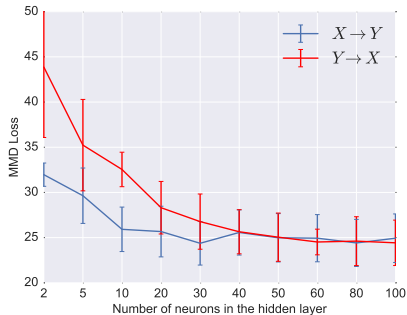
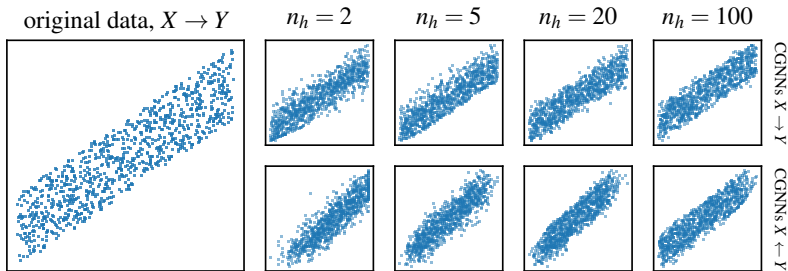
Modelling FCM with generative neural network

- Idea: approximate the continuous mechanism f_1, \dots, f_d with a set of one hidden layer neural networks $\hat{f} = (\hat{f}_1, \dots, \hat{f}_d)$
- Estimate FCMs C as $\hat{C} = (\hat{\mathcal{G}}, \hat{f}, \hat{Q})$:

$$\hat{X}_i \leftarrow \hat{f}_i(\hat{X}_{\text{Pa}(i; \hat{\mathcal{G}})}, \hat{E}_i), \hat{E}_i \sim \hat{Q}, \quad (1)$$

- We can draw a sample $\hat{x} = (\hat{x}_1, \dots, \hat{x}_d)$ from the distribution $\hat{P} := \hat{P}(X)$:
 1. Draw $\hat{e}_i \sim \hat{Q}$ for all $i = 1, \dots, d$.
 2. Construct $\hat{x}_i = \hat{f}_i(\hat{x}_{\text{Pa}(i; \hat{\mathcal{G}})}, \hat{e}_i)$ in the topological order of $\hat{\mathcal{G}}$

Complexity/reproduction trade-off



How to deal with the complexity/reproduction trade-off ?

Two ideas:

1. *Wrapper* approach

- CGNN - <https://arxiv.org/abs/1711.08936> (XCVML 2018)
- Limit the number of available hidden units n_h
- Explore the space of possible graph and find the DAG minimizing a MMD score

2. *Embedded* approach

- SAM: Structural Agnostic Model
- Use regularization to enforce automatically the sparsity of the graph and to choose the orientation of the edges
- Use GAN as a score to reproduce the distributions

SAM: Redefining causal graph recovery as a single optimization problem

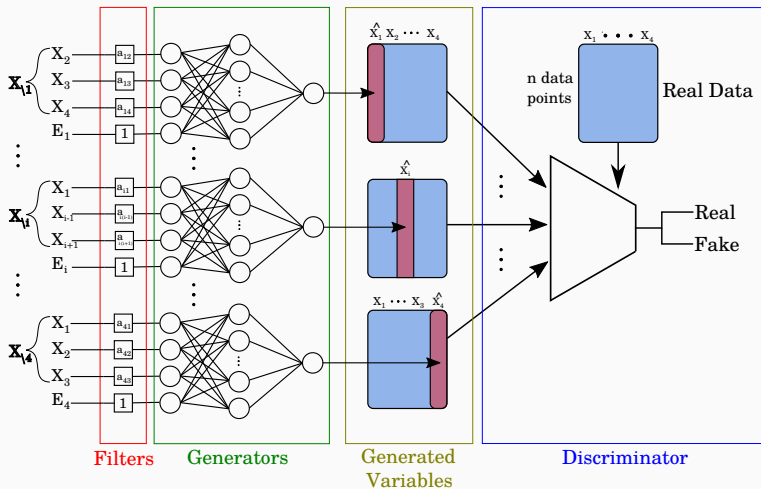
→ **Minimizing both structural and functional complexity of the graph**

Compared to score-based methods:

- For traditional score based methods, the graph structure search in the graph space (2^d given a skeleton)
- Retraining fully a regression model for each candidate graph

$$X_i = f_i(X_{\setminus i}, E_i), \quad (2)$$

Model Diagram



$$\hat{X}_i = m_i^\top \tanh \left(\bar{W}_i^\top (a_i \odot X) + n_i E_i + b_i \right) + \beta_i,$$

The generator i loss becomes :

$$\mathcal{L}_i = \frac{1}{2} \mathbb{E}_{x \sim p(\mathbf{x}_i)} [\log(\mathcal{D}(x|\mathbf{x}_{\setminus i}))] - \frac{1}{2} \mathbb{E}_{z \sim p(E_i)} [\log(1 - \mathcal{D}(\mathcal{G}(z|\mathbf{x}_{\setminus i})))] \quad (3)$$

Thus the total loss on generators:

$$\mathcal{L} = \underbrace{\sum_i^d \mathcal{L}_i}_{\text{adversarial generation loss}} + \underbrace{\sum_{i,j, i \neq j}^d |a_{ij}|}_{L_1 \text{ regularization}} \quad (4)$$

Output : adjacency matrix

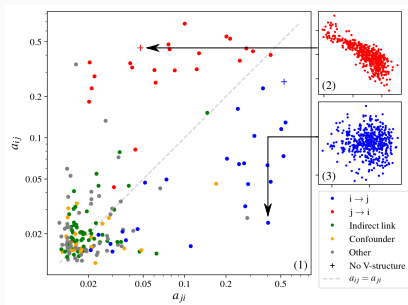
The matrix of a_{ij} represents the adjacency matrix of the graph.

$$\mathbf{A} = \underbrace{\left(\begin{array}{ccccc} 0 & a_{12} & a_{13} & \cdots & a_{1d} \\ \vdots & \ddots & & & \\ a_{j1} & \cdots & 0 & \cdots & a_{jd} \\ \vdots & & & \ddots & \\ a_{d1} & a_{d2} & \cdots & a_{d(d-1)} & 0 \end{array} \right)}_{d \text{ columns}} \left. \vphantom{\begin{pmatrix} 0 \\ \vdots \\ a_{j1} \\ \vdots \\ a_{d1} \end{pmatrix}} \right\} d \text{ rows}$$

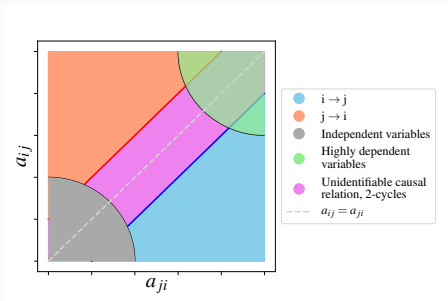
$$\boxed{X_i \rightarrow X_j \text{ if } a_{ij} < a_{ji}}$$

Experiments

Interpretation of A



(a) Projections of a_{ij}



(b) Causality map

Datasets

1. *Linear*: $X_i = \sum_{j \in \text{Pa}(i)} a_{i,j} X_j + E_i$.
2. *Sigmoid AM*: $X_i = \sum_{j \in \text{Pa}(i)} f_{i,j}(X_j) + E_i$, where $f_{i,j}(x_j) = a \cdot \frac{b \cdot (x_j + c)}{1 + |b \cdot (x_j + c)|}$ with $a \sim \text{Exp}(4) + 1$, $b \sim \mathcal{U}([-2, -0.5] \cup [0.5, 2])$ and $c \sim \mathcal{U}([-2, 2])$.
3. *Sigmoid Mix*: $X_i = f_i(\sum_{j \in \text{Pa}(i)} X_j + E_i)$, where f_i is as in the previous bullet-point.
4. *GP AM*: $X_i = \sum_{j \in \text{Pa}(i)} f_{i,j}(X_j) + E_i$ where $f_{i,j}$ is an univariate Gaussian process with a Gaussian kernel of unit bandwidth.
5. *GP Mix*: $X_i = f_i([X_{\text{Pa}(i)}, E_i])$, where f_i is a multivariate Gaussian process with a Gaussian kernel of unit bandwidth.
6. *Polynomial*: $X_i = \sum_{j \in \text{Pa}(i)} f_{i,j}(X_j) + E_i$, or $X_i = \sum_{j \in \text{Pa}(i)} f_{i,j}(X_j) \cdot E_i$, where $f_{i,j}$ is a random polynomial with a random degree in $[1, 4]$.
7. *NN*: $X_i = f_i(X_{\text{Pa}(i)}, E_i)$, where f_i a random single-hidden-layer neural network with 20 ReLU units.

Artificial graphs

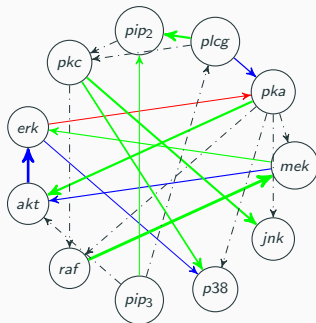
Table 1: Mean Average Precision (AP) on acyclic graphs. AP is between 0 and 1, 1 is best. Bold denotes best performance on the graph size. Underline denotes statistical significance at $p = 10^{-2}$. 55C with the HSIC test has not been evaluated on the graphs of 100 variables, being too computationally expensive.

Graph size	PC Gauss		PC HSIC		GES		MMHC		DAGL1		LINGAM		CAM		SAM	
	20	100	20	100	20	100	20	100	20	100	20	100	20	100	20	100
Linear	0.36	0.39	0.29	/	0.40	0.41	0.36	0.40	0.30	0.28	0.31	0.15	0.29	0.36	0.49	0.46
Sigmoid AM	0.28	0.24	0.33	/	0.18	0.12	0.31	0.28	0.19	0.15	0.19	0.09	0.72	<u>0.69</u>	0.73	0.57
Sigmoid Mix	0.22	0.21	0.25	/	0.21	0.15	0.22	0.23	0.16	0.17	0.12	0.05	0.15	0.19	<u>0.52</u>	<u>0.51</u>
GP AM	0.21	0.17	0.35	/	0.19	0.08	0.21	0.18	0.15	0.12	0.17	0.08	<u>0.96</u>	<u>0.95</u>	0.74	0.69
GP Mix	0.22	0.17	0.34	/	0.18	0.09	0.22	0.17	0.19	0.12	0.14	0.05	0.61	0.60	0.66	<u>0.66</u>
Polynomial	0.27	0.27	0.31	/	0.20	0.14	0.11	0.03	0.26	0.29	0.32	0.13	0.47	0.55	<u>0.65</u>	0.56
NN	0.40	0.40	0.38	/	0.42	0.37	0.11	0.03	0.43	0.50	0.36	0.19	0.22	0.32	<u>0.60</u>	0.55
Execution time	1s	23s	10h	>50h	<1s	5s	<1s	4s	2s	40s	2s	30s	2.5h	13h	1.2h	4h

Precision/recall score - causal protein network

Table 2: Mean Average Precision (AP) Cyto

	CCD	PC Gauss	GES	MMHC	DAGL1	LINGAM	CAM	CAM
Cyto	0.21	0.16	0.14	0.20	0.22	0.16	0.28	0.31



Towards improving SAM

What does represent $a_{i,j}$?

Confusion between:

1. The impact/amplitude of the causal effect
2. The existence of a causal relationship

Example :

$$X = 0.1 * Y + \mathcal{N}_{0,1}$$

Solution: Gumbel Softmax on edges

- By introducing a probability on whether the edge will be present or not, instead of the continuous a_{ij} .
- At each epoch a new graph $\hat{\mathcal{G}}$, where each edge (i, j) is drawn according to Bernoulli distribution with parameter p_{ij}
- These lambdas will be either 1 or 0 depending on the sampling made on a learned parameter p_{ij} through backpropagation [Jang et al., 2016]

$$a_{i,j} = \text{one hot}(\text{argmax}(g_0 + \log(1 - p_{ij}), g_1 + \log p_{ij})) \quad (5)$$

A skeleton recovery phase for better scaling (on big datasets)

Like many algorithms, the addition of a pruning phase before the algorithms, improves the consistency, computational efficiency and robustness of the results.

Adding such phase on SAM would help it to scale better (currently d^2)

An (optional) DAG constraint

If was generated by a DAG (Directed Acyclic Graph), an additional term can be added to recover a DAG [Zheng et al., 2018]:

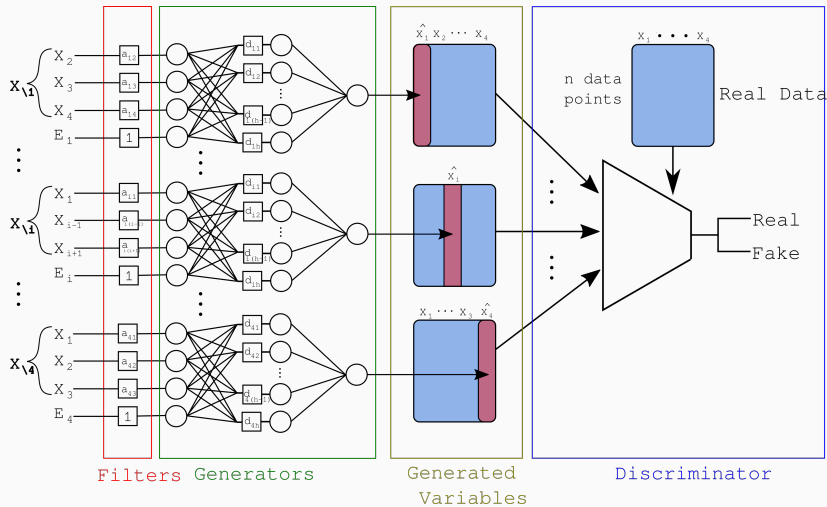
$$\mathcal{L}_{DAG} = \sum_{k=1}^d \frac{\text{tr}(\mathbf{A}^k)}{k!}$$

An automatic regularization of the functional complexity

Before, the number of hidden units was fixed and dependent of the mechanism complexity that generated the data.

Adding a sampling of units, like the a coefficients, adds a control term for this functional complexity.

$$\hat{X}_j = m_j^\top \tanh \left(\bar{W}_j^\top (a_j \odot X) + n_j E_j + b_j \right) \odot d_j + \beta_j, \quad (6)$$



What is usable right now

All the presented framework is available on GitHub at :

<https://github.com/Diviyan-Kalainathan/CausalDiscoveryToolbox>

It includes multiple algorithms as well as tools for graph structure.

Questions?



Hoyer, P. O., Janzing, D., Mooij, J. M., Peters, J., and Schölkopf, B. (2009).

Nonlinear causal discovery with additive noise models.

In *Advances in neural information processing systems*, pages 689–696.



Jang, E., Gu, S., and Poole, B. (2016).

Categorical reparameterization with gumbel-softmax.

arXiv preprint arXiv:1611.01144.



Spirtes, P., Glymour, C. N., and Scheines, R. (2000).

Causation, prediction, and search.

MIT press.



Tsamardinos, I., Brown, L. E., and Aliferis, C. F. (2006).

The max-min hill-climbing bayesian network structure learning algorithm.

Machine learning, 65(1):31–78.



Zheng, X., Aragam, B., Ravikumar, P., and Xing, E. P.
(2018).

**DAGs with NO TEARS: Smooth Optimization for
Structure Learning.**

ArXiv e-prints.