

Causal Modeling

Michèle Sebag

TAU: Tackling the Underspecified

CNRS – INRIA – LRI – Université Paris-Saclay



université
PARIS-SACLAY

Paris – Jan. 13th, 2020

A Case of Irrational Scientific Exuberance

- ▶ Underspecified goals Big Data cures everything
- ▶ Underspecified limitations Big Data can do anything (if big enough)
- ▶ Underspecified caveats Big Data and Big Brother

Wanted: An AI with common decency

- ▶ Fair no biases
- ▶ Accountable models can be explained
- ▶ Transparent decisions can be explained
- ▶ Robust w.r.t. malicious examples

ML & AI, 2

In practice

- ▶ Data are ridden with biases
- ▶ Learned models are biased (prejudices are transmissible to AI agents)
- ▶ Issues with robustness
- ▶ Models are used out of their scope

More

- ▶ C. O'Neill, *Weapons of Math Destruction*, 2016
- ▶ Zeynep Tufekci, *We're building a dystopia just to make people click on ads*, Ted Talks, Oct 2017.

Machine Learning: discriminative or generative modelling

Given a training set

iid samples $\sim P(X, Y)$

$$\mathcal{E} = \{(\mathbf{x}_i, y_i), \mathbf{x}_i \in \mathbb{R}^d, i \in [[1, n]]\}$$

Find

- ▶ Supervised learning: $\hat{h} : X \mapsto Y$ or $\hat{P}(Y|X)$
- ▶ Generative model $\hat{P}(X, Y)$

Predictive modelling might be based on correlations

If umbrellas in the street, Then it rains



The implicit big data promise:

If you can predict what will happen,
then how to make it happen what you want ?

Knowledge → **Prediction** → **Control**

ML models will be expected to support *interventions*:

- ▶ health and nutrition
- ▶ education
- ▶ economics/management
- ▶ climate

The implicit big data promise, 2

Intervention

Pearl 2009

Intervention $do(X = a)$ forces variable X to value a

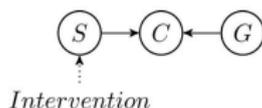
Direct cause $X \rightarrow Y$

$$P_{Y|do(X=a, Z=c)} \neq P_{Y|do(X=b, Z=c)}$$

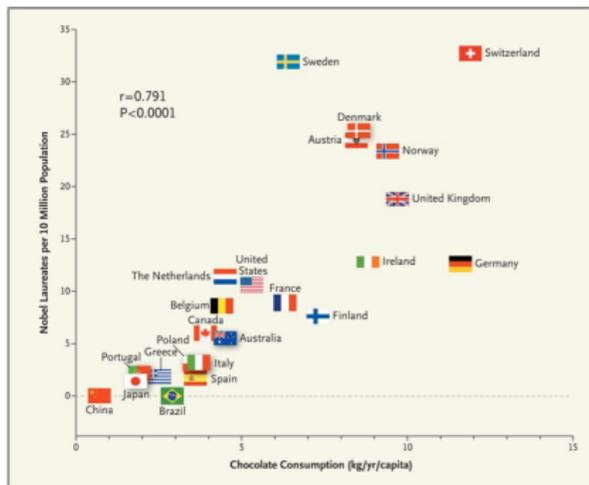
Example

C: Cancer, S: Smoking, G: Genetic factors

$$P(C|do\{S = 0, G = 0\}) \neq P(C|do\{S = 1, G = 0\})$$



Correlations do not support interventions



F. H. Messerli: *Chocolate Consumption, Cognitive Function, and Nobel Laureates*, N Engl J Med 2012

Causal models are needed to support interventions

*Consumption of chocolate enables to predict # of Nobel prizes
but eating more chocolates does not increase # of Nobel prizes*

Predictive model \nrightarrow Causal model

Consider

$$X, E_Y, E_Z \sim \text{Uniform}(0, 1),$$

$$Y \leftarrow 0.5X + E_Y,$$

$$Z \leftarrow Y + E_Z,$$

with $E_Y, E_Z \sim \mathcal{N}(0, 1)$ (noise)

Predicting Y

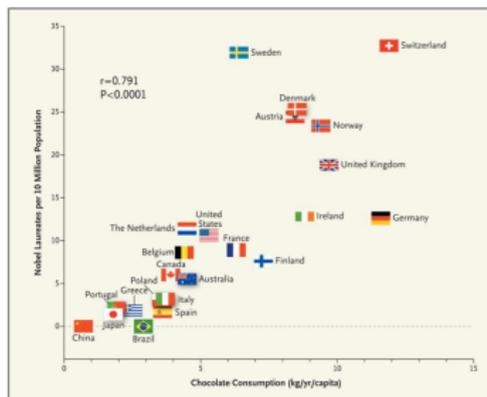
$$\hat{Y} = 0.25X + 0.5Z$$

If interpreted as a causal model, suggests that Y depends on Z .

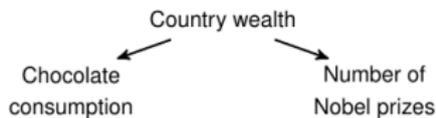
Issue

Causes can often be predicted from their effects

Confounders: When correlations do not imply causality



F. H. Messeri: Chocolate Consumption, Cognitive Function, and Nobel Laureates, N Engl J Med 2012



Tentative explanation

- ▶ Both effects of a same cause, $C \not\perp N$.
- ▶ But C and N are conditionally independent given W

$$C \perp\!\!\!\perp N | W$$

Causality and paradoxes

Facts

- ▶ If mother smokes, child weight tends to be low
- ▶ If child weight is low, more health problems
- ▶ However, low child weight AND mother smokes $>$ low child weight

Interpretation mother smoking beneficial to child's health ?

Explaining away

Many possible causes for low child weight

Many of these severely affect child's health (genetic diseases)

Compared to these, mother smoking is rather a good news...

An AI with common decency

Desired properties

- ▶ Fair
- ▶ Accountable
- ▶ Transparent
- ▶ Robust

no biases

models can be explained

decisions can be explained

w.r.t. malicious examples

Relevance of Causal Modeling

- ▶ Decreased sensitivity wrt data distribution
- ▶ Support interventions
- ▶ Hopes of explanations / bias detection

clamping variable value

1.State of the art

Causal Modelling

The Causal Discovery Setting

Assume random variables

X_1, \dots, X_d : random variables

and a sample of their joint distribution

$$\mathcal{D} = \{\mathbf{x}_i, i = 1 \dots n\}$$

to be given.

Formal background: Overview

1. Key concepts
2. Framework
3. Approaches

Key concepts: 1. Dependence among pairs of variables

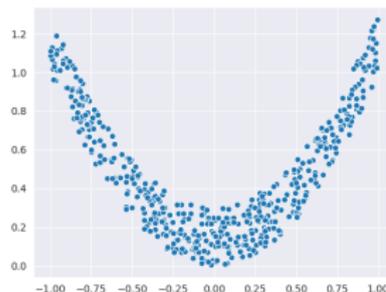
Independent variables X and Y ($X \perp\!\!\!\perp Y$)

$$X \perp\!\!\!\perp Y \text{ iff } P(X, Y) = P(X) \cdot P(Y)$$

Dependency tests

- Correlation

limited to linear dependencies



$$Y = X^2 + E$$
$$\text{Correlation}(X, Y) \approx 0$$

Key concepts: 1. Dependence among pairs of variables

Independent variables X and Y ($X \perp\!\!\!\perp Y$)

$$X \perp\!\!\!\perp Y \text{ iff } P(X, Y) = P(X).P(Y)$$

Dependency tests

- Correlation limited to linear dependencies
- HSIC, Hilbert-Schmitt Independence Criterion [Gretton et al., 2005]

$$HSIC(P_{XY}, \mathcal{F}, \mathcal{G}) := \|C_{XY}\|^2$$

where $\|\cdot\|$ denotes the Hilbert-Schmidt norm, and C_{XY} a kernel based covariance operator and \mathcal{F}, \mathcal{G} two RKHSs.

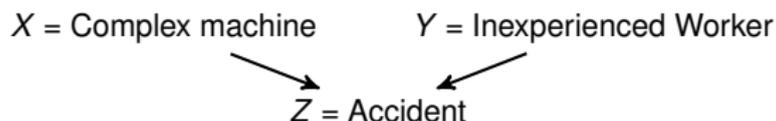
Key concepts: 2. Conditional Dependence/Independence

Conditional independence a.k.a. hidden confounder

Key concepts: 2. Conditional Dependence/Independence

Conditional independence a.k.a. hidden confounder

Conditional dependence a.k.a. V-structure



X and Y are independent; but given $Z = true$ they are not independent (either the machine is complex or the worker is inexperienced...)

Definition of causal relationship

Definition of intervention

$do(X = 1)$ forces variable X to value 1

[Pearl, 2009]

Definition of causal relationship

X is a direct cause of Y ($X \rightarrow Y$) iff
all other variables Z being constant,

$$P_{Y|do(X=1, \dots, Z=c)} \neq P_{Y|do(X=0, \dots, Z=c)} \quad (1)$$

Definition of causal relationship

Definition of intervention

$do(X = 1)$ forces variable X to value 1

[Pearl, 2009]

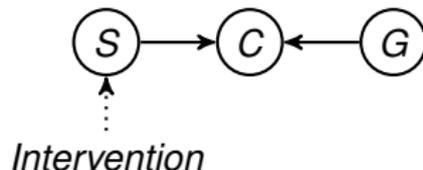
Definition of causal relationship

X is a direct cause of Y ($X \rightarrow Y$) iff
all other variables Z being constant,

$$P_Y|do(X=1, \dots, Z=c) \neq P_Y|do(X=0, \dots, Z=c) \quad (1)$$

Example C : Cancer, S : Smoking, G : Genetic factors.

$$P(C|do\{S = 0\}, G) \neq P(C|do\{S = 1\}, G)$$



Markov equivalence class and V-structure

Markov Equivalent Class: $A \perp\!\!\!\perp C \mid B$ and $A \not\perp\!\!\!\perp C$

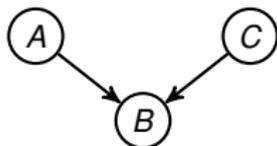


Markov equivalence class and V-structure

Markov Equivalent Class: $A \perp\!\!\!\perp C \mid B$ and $A \not\perp\!\!\!\perp C$



V-Structure: $A \not\perp\!\!\!\perp C \mid B$ and $A \perp\!\!\!\perp C$



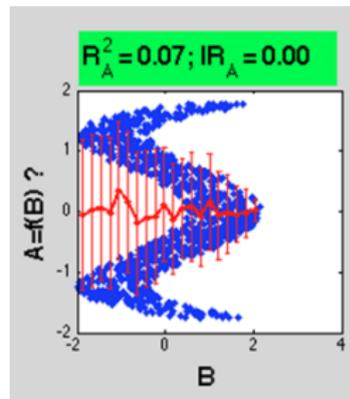
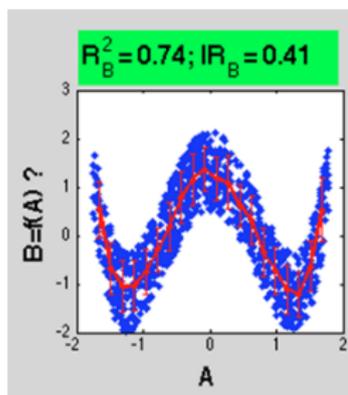
[Spirtes et al., 2000, Spirtes and Zhang, 2016]

Key concepts: 3. Causality with distributional asymmetry

Leveraging Occam's razor principle;

[Janzing, 2019]

→ the causal model as the one being the simplest model that fits the data.



Framework: Functional Causal Models (FCMs)

Given X_1, \dots, X_d ,

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

with $X_{\text{Pa}(i; \mathcal{G})}$ the set of parents of X_i in \mathcal{G} (= causes of X_i),

E_i a random independent noise variable modeling the unobserved other causes,

f_i a deterministic function: the causal mechanism

Framework: Functional Causal Models (FCMs)

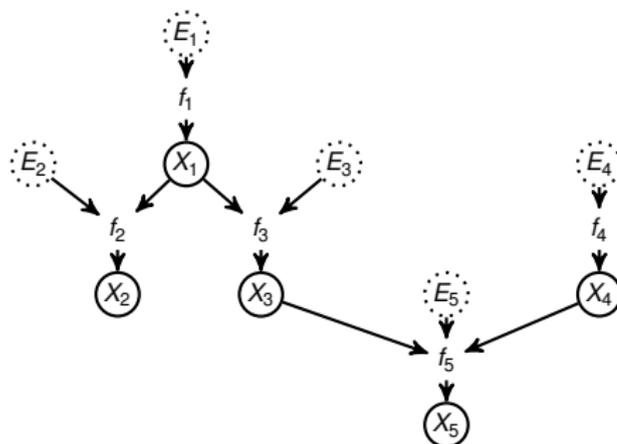
Given X_1, \dots, X_d ,

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

with $X_{\text{Pa}(i; \mathcal{G})}$ the set of parents of X_i in \mathcal{G} (= causes of X_i),

E_i a random independent noise variable modeling the unobserved other causes,

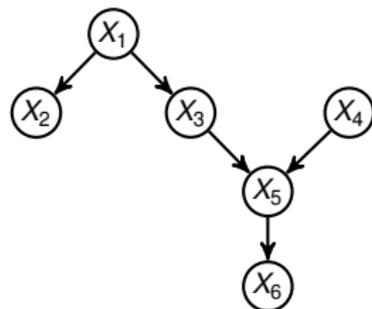
f_i a deterministic function: the causal mechanism



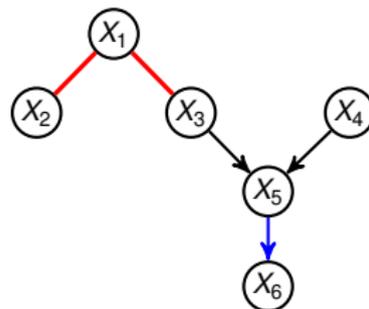
$$\begin{cases} X_1 = f_1(E_1) \\ X_2 = f_2(X_1, E_2) \\ X_3 = f_3(X_1, E_3) \\ X_4 = f_4(E_4) \\ X_5 = f_5(X_3, X_4, E_5) \end{cases}$$

Key approach 1: Constraint-based methods

Constraint-based methods, through V-Structures and constraint propagation, output a **CPDAG** (Completed Partially Directed Acyclic Graph).



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



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

Ex: Peter-Clark Algorithm (PC)

[Spirtes et al., 2000]

Non-linear extensions (CI tests): PC-HSIC (KCI-test), PC-RCIT

[Zhang et al., 2012, Strobl et al., 2017]

Key approach 2: Score-based methods

Objective function to optimize such as the Bayesian Information Criterion (BIC):

$$BIC(\mathcal{G}) = -2 \ln L + k * \ln n$$

with L : Likelihood of the model, k : number of parameters, n : Number of samples

Key approach 2: Score-based methods

Objective function to optimize such as the Bayesian Information Criterion (BIC):

$$BIC(\mathcal{G}) = -2 \ln L + k * \ln n$$

with L : Likelihood of the model, k : number of parameters, n : Number of samples

The graph is optimized with the operators:

- add edge
- remove edge
- revert edge

Ex: Greedy Equivalence Search (GES)

[Chickering, 2002]

Limitations

- Computational cost dependent on the type of test/scoring method used
- Data hungry
- Identifiability issues

Limitations

- Computational cost dependent on the type of test/scoring method used
- Data hungry
- Identifiability issues

Example

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

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

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



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

Key approach 3: Global optimization

Assuming linear causal mechanisms, the causal mechanisms can be formulated in terms of linear algebra.

$$\mathbf{X} = B^T \mathbf{X} + E$$

And estimate the B matrix, through ICA for LiNGAM

[Shimizu et al., 2006, Hyvärinen and Pajunen, 1999]

→ Graphical models

[Pearl, 2009, Friedman et al., 2008]

Ex: Max-Min Hill-Climbing (MMHC)

[Tsamardinos et al., 2006]

Concave penalized Coordinate Descent (CCDr)

[Aragam and Zhou, 2015]

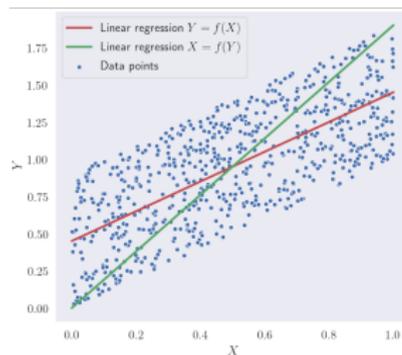
Key approach 4: Exploiting asymmetries in the distribution

→ If no v-structure available or causal discovery with 2 variables: leverage asymmetries in the distributions.

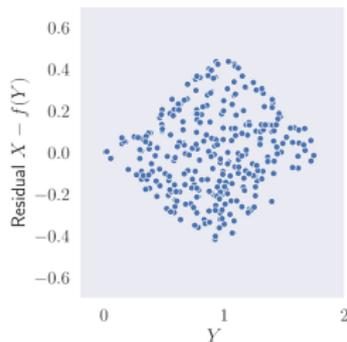
Additive noise model (ANM):

[Hoyer et al., 2009]

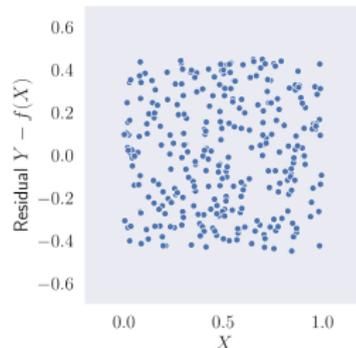
$$Y = f(X) + E$$



Original data



Residuals of X=g(Y)



Residuals of Y=f(X)

Ex: Post Non-Linear model (PNL), GPI

[Zhang and Hyvärinen, 2010, Stegle et al., 2010]

Limitations of asymmetry-based approaches

- Restrictive assumptions on the type of causal mechanisms
- Does not take into account conditional independence relations.

[Zhang and Hyvärinen, 2009]

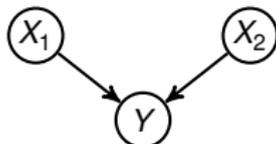
Limitations of asymmetry-based approaches

- Restrictive assumptions on the type of causal mechanisms
- Does not take into account conditional independence relations.

[Zhang and Hyvärinen, 2009]

Example

$$X_1, X_2, E_{X_1} \sim \text{Gaussian}(0, 1), X_1 \perp\!\!\!\perp E_{X_1}, X_2 \perp\!\!\!\perp E_{X_1}$$
$$Y \leftarrow 0.5X_1 + X_2 + E_{X_1}$$



(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

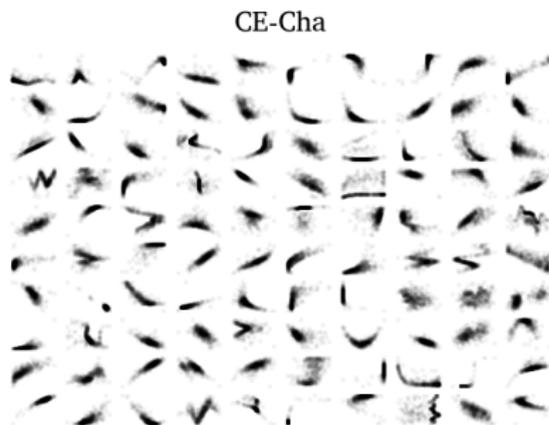
Key approach 5: Supervised learning for causation identification

Reformulate the pairwise cause-effect problem as a pattern recognition problem:

[Guyon, 2013, Guyon, 2014]

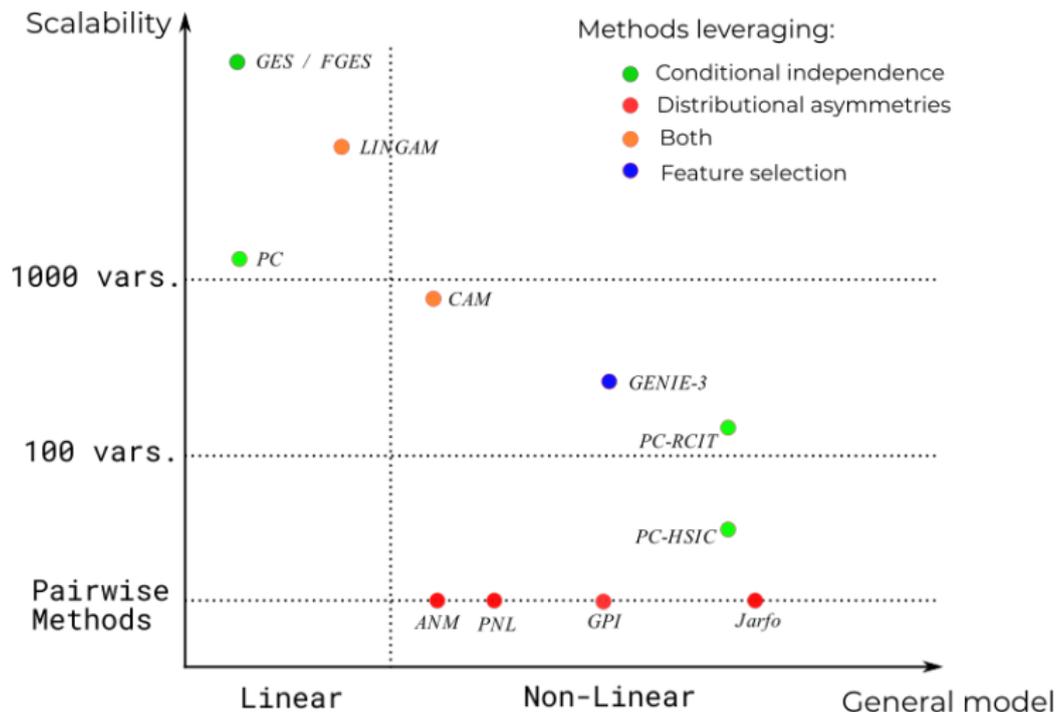
Given a pair of variables (X, Y) :

Label: $X \rightarrow Y$ or $Y \rightarrow X$ or $X \leftrightarrow Y$

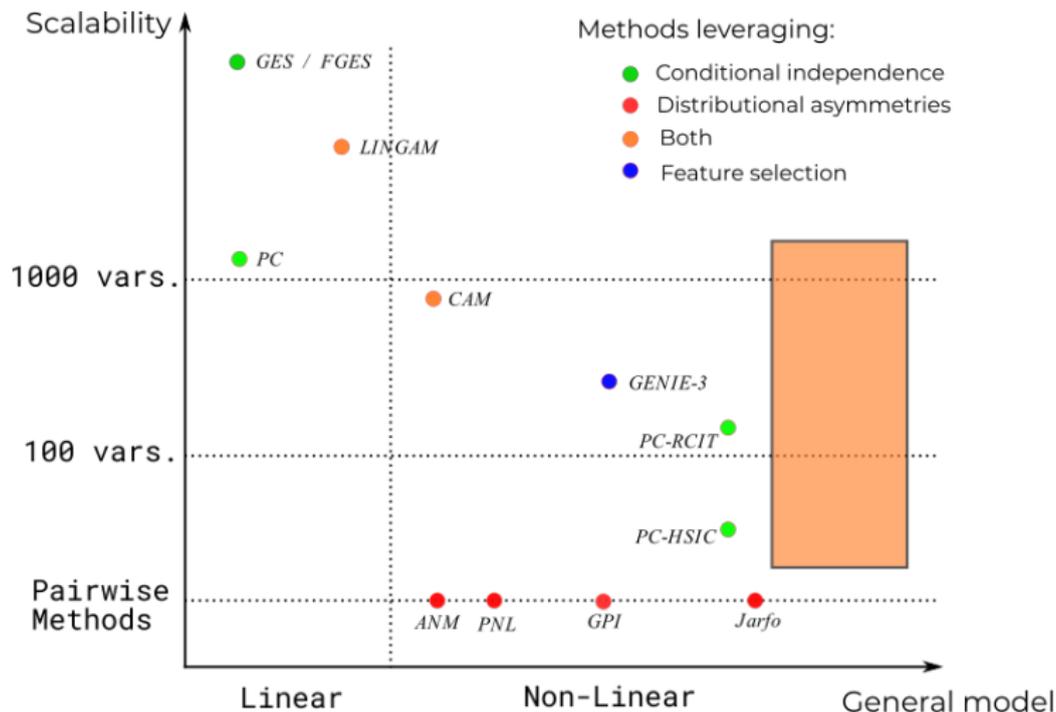


Example pairs of the cause-effect challenge

State of the art: summary



State of the art: summary



Motivation

State of the art

Formal Background

The cause-effect pair challenge

The general setting

Causal Generative Neural Nets

Applications

Human Resources

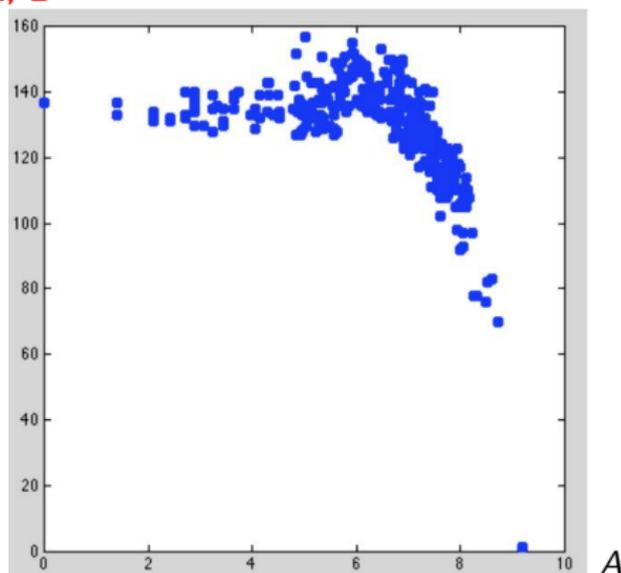
Food and Health

Discussion

Causality: What ML can bring ?

Each point: sample of the joint distribution $P(A, B)$.

Given variables A, B



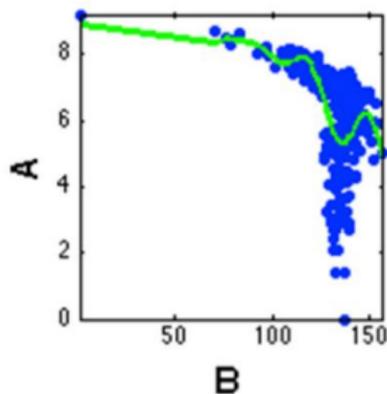
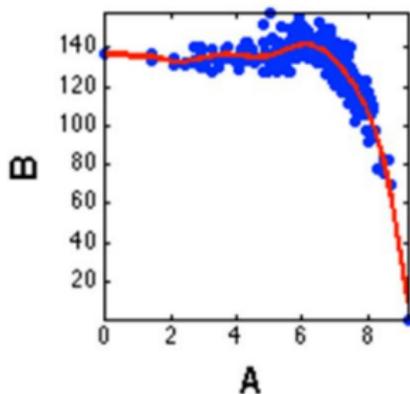
Causality: What ML can bring, follow'd

Given A , B , consider models

- ▶ $A = f(B)$
- ▶ $B = g(A)$

Compare the models

Select the best model: $A \rightarrow B$



Causality: What ML can bring, follow'd

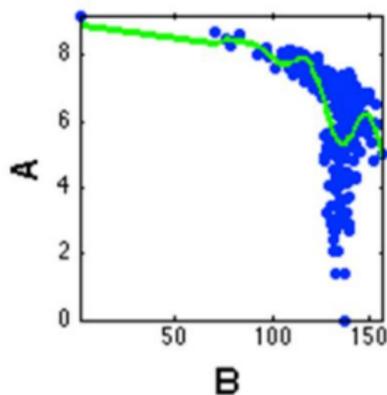
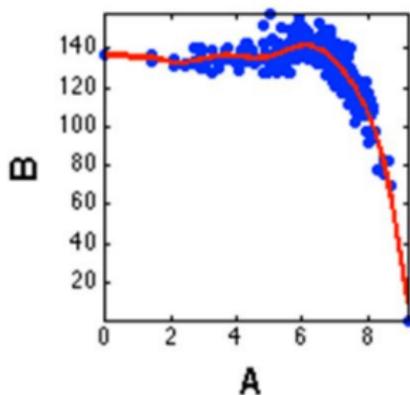
Given A , B , consider models

▶ $A = f(B)$

▶ $B = g(A)$

Compare the models

Select the best model: $A \rightarrow B$



A : Altitude, B : Temperature

Each point = (altitude, average temperature of a city)

Causality: A machine learning-based approach

Guyon et al, 2014-2015

Pair Cause-Effect Challenges

- ▶ Gather data: a sample is a pair of variables (A_i, B_i)
- ▶ Its label ℓ_i is the “true” causal relation (e.g., age “causes” salary)

Input

$$\mathcal{E} = \{(A_i, B_i, \ell_i), \ell_i \text{ in } \{\rightarrow, \leftarrow, \perp\perp\}\}$$

Example A_i, B_i	Label ℓ_i
A_i causes B_i	\rightarrow
B_i causes A_i	\leftarrow
A_i and B_i are independent	$\perp\perp$

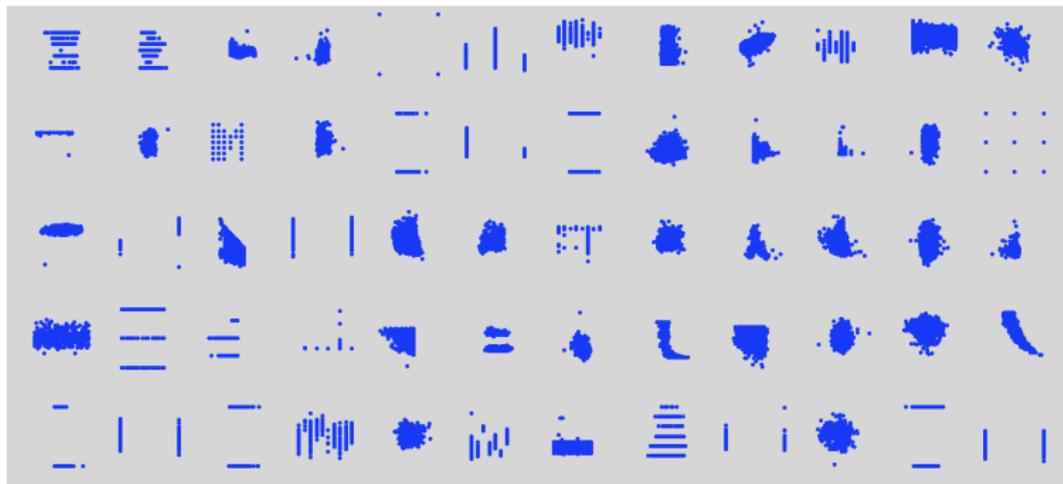
Output

using supervised Machine Learning

Hypothesis : $(A, B) \mapsto \text{Label}$

Causality: A machine learning-based approach, 2

Guyon et al, 2014-2015

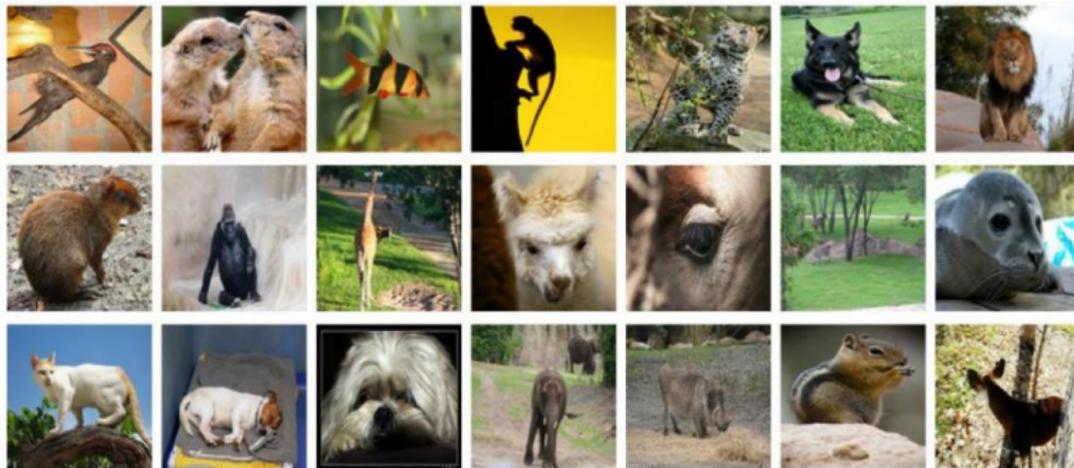


The Cause-Effect Pair Challenge

Learn a **causality classifier** (causation estimation)

- ▶ Like for any supervised ML problem from images

ImageNet 2012



More

- ▶ Guyon et al., eds, *Cause Effect Pairs in Machine Learning*, 2019.

Motivation

State of the art

Formal Background

The cause-effect pair challenge

The general setting

Causal Generative Neural Nets

Applications

Human Resources

Food and Health

Discussion

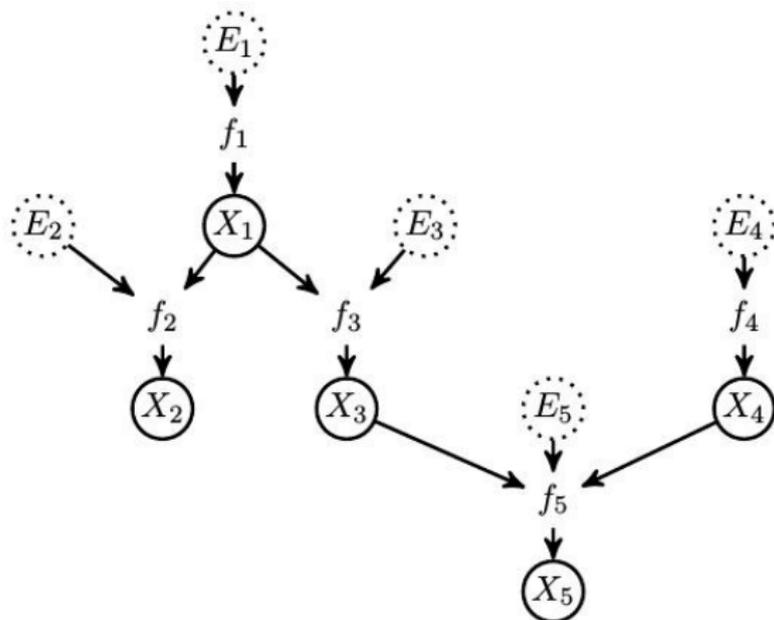
Functional Causal Models, a.k.a. Structural Equation Models

Pearl 00-09

$$X_i = f_i(\text{Pa}(X_i), E_i)$$

$\text{Pa}(X_i)$: Direct causes for X_i

E_i : noise variables, all unobserved influences



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

Tasks

- ▶ Finding the structure of the graph (no cycles)
- ▶ Finding functions (f_i)

Conducting a causal modelling study

Spirtes et al. 01; Tsamardinos et al., 06; Hoyer et al. 09
Daniusis et al., 12; Mooij et al. 16

Milestones

- ▶ Testing bivariate independence (statistical tests)
find edges
- ▶ Conditional independence
prune the edges
- ▶ Full causal graph modelling
orient the edges

$$X - Y; Y - Z$$

$$X \perp\!\!\!\perp Z | Y$$

$$X \rightarrow Y \rightarrow Z$$

Challenges

- ▶ Computational complexity
- ▶ Conditional independence: data hungry tests
- ▶ Assuming causal sufficiency

tractable approximation

can be relaxed

$X - Y$ independence

$$P(X, Y) \stackrel{?}{=} P(X).P(Y)$$

Categorical variables

- ▶ Entropy $H(X) = -\sum_x p(x)\log(p(x))$
x: value taken by X , $p(x)$ its frequency
- ▶ Mutual information $M(X, Y) = H(X) + H(Y) - H(X, Y)$
- ▶ Others: χ^2 , G-test

Continuous variables

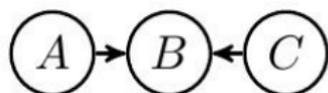
- ▶ t-test, z-test
- ▶ Hilbert-Schmidt Independence Criterion (HSIC) Gretton et al., 05

$$\text{Cov}(f, g) = \mathbb{E}_{x,y}[f(x)g(y)] - \mathbb{E}_x[f(x)]\mathbb{E}_y[g(y)]$$

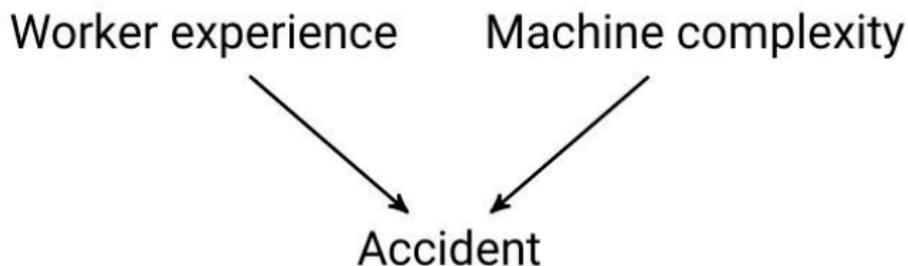
- ▶ Given $f : X \mapsto \mathbb{R}$ and $g : Y \mapsto \mathbb{R}$
- ▶ $\text{Cov}(f, g) = 0$ for all f, g iff X and Y are independent

Find V-structure: $A \perp\!\!\!\perp C$ and $A \not\perp\!\!\!\perp C|B$

Explaining away causes



Example



Motivation

State of the art

Formal Background

The cause-effect pair challenge

The general setting

Causal Generative Neural Nets

Applications

Human Resources

Food and Health

Discussion

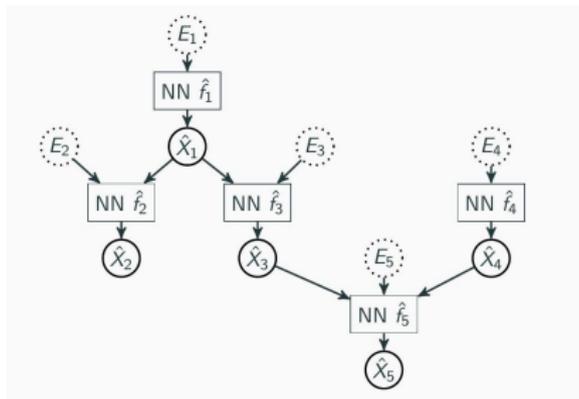
Causal Generative Neural Network

Goudet et al. 17

Principle

- ▶ Given skeleton
- ▶ Given X_i and candidate $Pa(i)$
- ▶ Learn $f_i(Pa(X_i), E_i)$ as a generative neural net
- ▶ Train and compare candidates based on scores

given or extracted



NB

- ▶ Can handle confounders (X_1 missing $\rightarrow (E_2, E_3 \rightarrow E_{2,3})$)

Causal Generative Neural Network (2)

Training loss

- ▶ Observational data $\mathbf{x} = \{[x_1, \dots, x_n]\}$ x_i in $\mathbb{R}^{* * d}$
- ▶ (Graph, \hat{f}) $\hat{\mathbf{x}} = \{[\hat{x}_1, \dots, \hat{x}_{n'}]\}$ \hat{x}_i in $\mathbb{R}^{* * d}$
- ▶ Loss: Maximum Mean Discrepancy ($\mathbf{x}, \hat{\mathbf{x}}$) (+ parsimony term), with k kernel (Gaussian, multi-bandwidth)

$$\text{MMD}_k(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{n^2} \sum_{i,j} k(x_i, x_j) + \frac{1}{n'^2} \sum_{i,j} k(\hat{x}_i, \hat{x}_j) - \frac{2}{n \times n'} \sum_{i=1}^n \sum_{j=1}^{n'} k(x_i, \hat{x}_j)$$

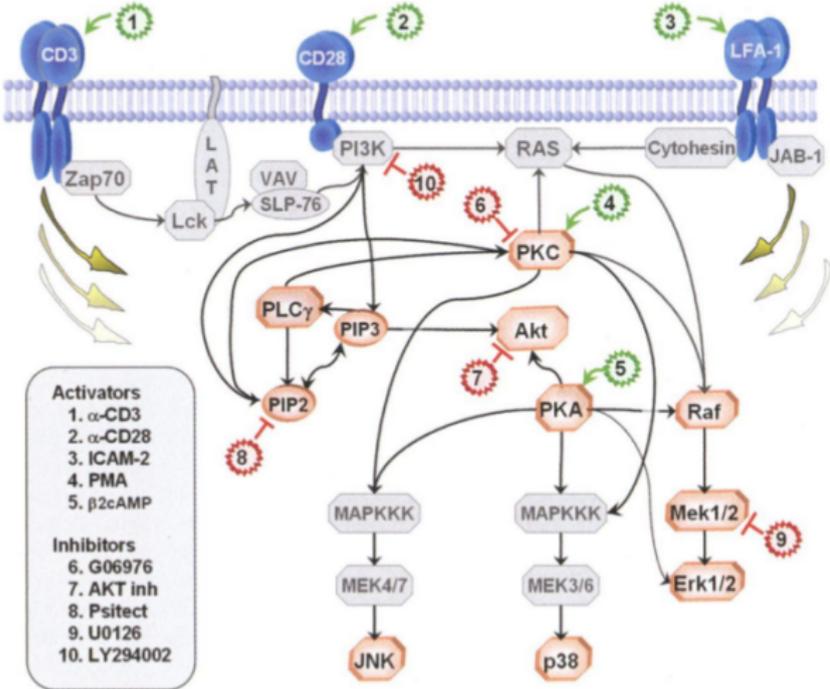
- ▶ For $n, n' \rightarrow \infty$

$$\text{MMD}_k(\mathbf{x}, \hat{\mathbf{x}}) = 0 \Rightarrow \mathcal{D}(\mathbf{x}) = \mathcal{D}(\hat{\mathbf{x}})$$

Gretton 07

Results on real data: causal protein network

Sachs et al. 05



Edge orientation task

All algorithms start from the skeleton of the graph

method	AUPR	SHD	SID
<i>Constraints</i>			
PC-Gauss	0.19 (0.07)	16.4 (1.3)	91.9 (12.3)
PC-HSIC	0.18 (0.01)	17.1 (1.1)	90.8 (2.6)
<i>Pairwise</i>			
ANM	0.34 (0.05)	8.6 (1.3)	85.9 (10.1)
Jarfo	0.33 (0.02)	10.2 (0.8)	92.2 (5.2)
<i>Score-based</i>			
GES	0.26 (0.01)	12.1 (0.3)	92.3 (5.4)
LiNGAM	0.29 (0.03)	10.5 (0.8)	83.1 (4.8)
CAM	0.37 (0.10)	8.5 (2.2)	78.1 (10.3)
CGNN ($\widehat{\text{MMD}}_k$)	0.74* (0.09)	4.3* (1.6)	46.6* (12.4)

AUPR: Area under the Precision Recall Curve

SHD: Structural Hamming Distance

SID: Structural intervention distance

Limitations

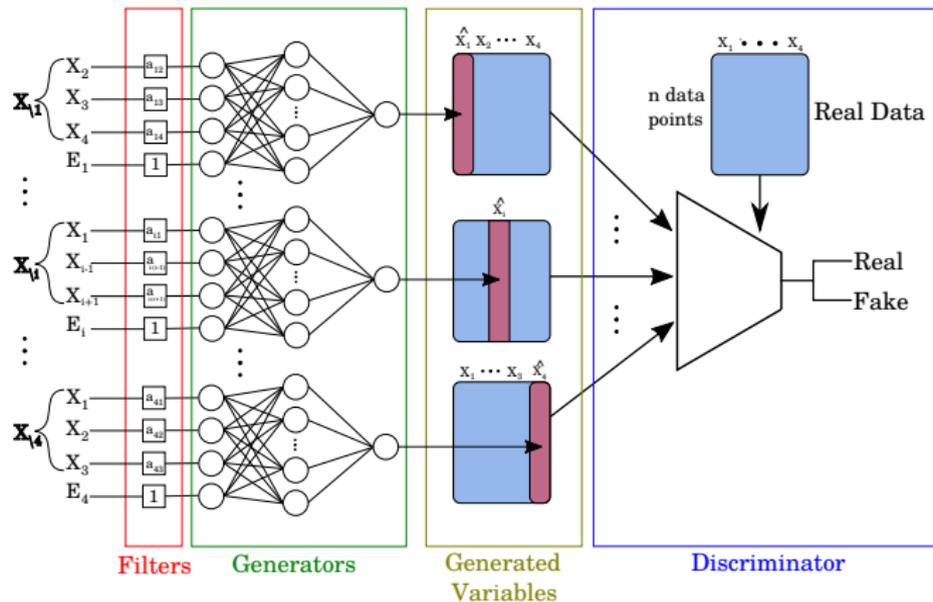
- ▶ Combinatorial search in the structure space
- ▶ Retraining fully the NN for each candidate graph
- ▶ MMD Loss is $O(n^2)$
- ▶ Limited to DAG

Structure Agnostic Modeling

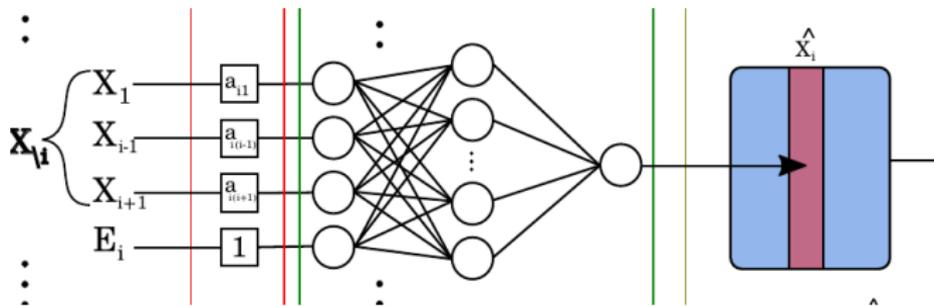
Kalainathan et al. 18

Goal: A generative model

- + Does not require CPDAG as input
- + Avoids combinatorial search for structure
- Less computationally demanding



Structure Agnostic Modeling, 2



The i -th neural net

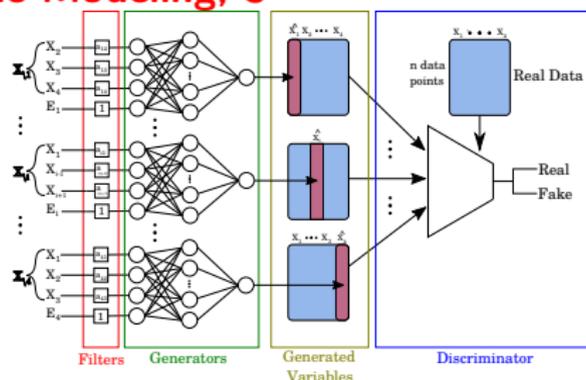
- ▶ Learns conditional distribution $P(X_i|X_{\setminus i})$ as $\hat{f}_i(X_{\setminus i}, E_i)$
- ▶ Filter variables $a_{i,j}$ are used to enforce sparsity (Lasso-like, next slide)
- ▶ 1st non-linear layer builds features $\phi_{i,k}$, 2nd layer builds linear combination of features:

$$f_i(X_{\setminus i}, E_i) = \sum \beta_{i,k} \phi_{i,k}(a_{i,1}X_1, \dots, a_{i,d}X_d, E_i)$$

In the large sample limit, $a_{i,j} = 1$ iff $X_j \in MB(X_j)$

Yu et al. 18

Structure Agnostic Modeling, 3



Given observational data $\{x_1, \dots, x_n\} \sim P(X_1, \dots, X_d)$

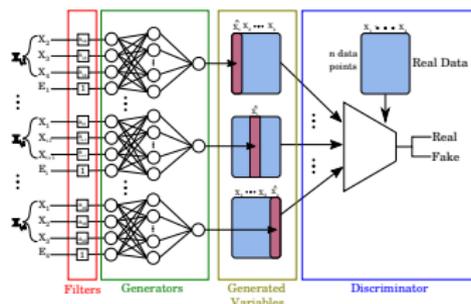
x_i in \mathbb{R}^d

Adversarial learning

- ▶ Generate $\{\tilde{x}_i^{(j)}\}$ with j -th component of $\tilde{x}_i^{(j)}$ set to $\hat{f}_i(x_i, \epsilon)$, $\epsilon \sim \mathcal{N}(0, 1)$
- ▶ Discriminator D among observational data $\{x_i\}$ and generated data $\{\tilde{x}_i^{(j)}, i = [[1, n]], j = [[1, d]]\}$
- ▶ Learning criterion (adversarial + sparsity)

$$\min \left(\text{Accuracy}(D) + \lambda \sum_{i,j} |a_{i,j}| \right)$$

Structure Agnostic Modeling, 4



Learning criterion $\min \left(\text{Accuracy}(D) + \lambda \sum_{i,j} |a_{i,j}| \right)$

Competition between discriminator and sparsity term $\sum \|\mathbf{a}\|_1$

- ▶ Avoids combinatorial search for structure
- ▶ Cycles are possible
- ▶ DAGness achieved by enforcing constraints on trace of $A = (a_{i,j})$ and A^k

Quantitative benchmark - artificial DAG

Directed **acyclic** artificial graphs (DAG) of 20 variables

	PC Gauss	PC HSIC	GES	MMHC	DAGL1	LINGAM	CAM	SAM
Linear	0.36	0.29	0.40	0.36	0.30	0.31	0.29	0.49
Sigmoid AM	0.28	0.33	0.18	0.31	0.19	0.19	0.72	0.73
Sigmoid Mix	0.22	0.25	0.21	0.22	0.16	0.12	0.15	<u>0.52</u>
GP AM	0.21	0.35	0.19	0.21	0.15	0.17	<u>0.96</u>	0.74
GP Mix	0.22	0.34	0.18	0.22	0.19	0.14	0.61	0.66
Polynomial	0.27	0.31	0.20	0.11	0.26	0.32	0.47	<u>0.65</u>
NN	0.40	0.38	0.42	0.11	0.43	0.36	0.22	<u>0.60</u>
Execution time	1s	10h	<1s	<1s	2s	2s	2.5h	1.2h

Quantitative benchmark - artificial DG (with cycles)

Directed **cyclic** artificial graphs of 20 variables

	CCD	PC Gauss	GES	MMHC	DAGL1	LINGAM	CAM	SAM
Linear	0.44	0.44	0.20	0.34	0.26	0.19	0.23	0.51
Sigmoid AM	0.31	0.31	0.16	0.32	0.17	0.24	0.37	0.47
Sigmoid Mix	0.31	0.35	0.18	0.34	0.19	0.17	0.22	0.49
GP AM	0.30	0.32	0.17	0.30	0.15	0.23	0.50	0.56
GP Mix	0.24	0.25	0.15	0.24	0.16	0.18	0.26	0.49
Polynomial	0.25	0.33	0.20	0.25	0.17	0.22	0.33	0.42
NN	0.25	0.18	0.18	0.24	0.18	0.16	0.22	0.40
Execution time	1s	1s	<1s	<1s	2s	2s	2.5h	1.2h

Motivation

State of the art

Formal Background

The cause-effect pair challenge

The general setting

Causal Generative Neural Nets

Applications

Human Resources

Food and Health

Discussion

Causal Modeling and Human Resources

Known:

- A Quality of life at work
- B Economic performance
- ▶ ... are correlated

employee's perspective

firm's perspective

Question: Are there causal relationships ?

$A \rightarrow B$; or $B \rightarrow A$; or $\exists C / C \rightarrow A$ and $C \rightarrow B$

Data

- ▶ Polls from Ministry of Labor
- ▶ Gathered by Group Alpha Secafi (trade union advisor)
- ▶ Tax files + social audits for 408 firms

Economic sectors: low tech, medium-low, medium-high and high-tech.

Variables

Economic indicators

- ▶ Total number of employees
- ▶ Capitalistic intensity, Total payroll, Gini index
- ▶ Average salary (of workers, technicians, managers)
- ▶ Productivity, Operating profits, Investment rate

People

- ▶ Average age, Average seniority, Physical effort,
- ▶ Permanent contract rate, Manager rate, Fixed-term contract rate, Temporary job rate, Shift and night work, Turn-over
- ▶ Vocational education effort, duration of stints, Average stint rate (for workers, technicians, managers);

Variables, cont'd

Quality of life at work

- ▶ Frequency & Gravity of work injuries, Safety expenses, Safety training expenses
- ▶ Absenteeism (diseases), Occupational-related diseases
- ▶ Resignation rate, Termination rate, Participation rate
- ▶ Subsidy to the works council

Men/Women

- ▶ Percentage of women (employees, managers)
- ▶ Wage gap between women and men (average, for workers, technicians, managers)

General Causal Relations

Access to training ↗

- ▶ ↘ Gravity of work injuries
- ▶ ↘ Occupational-related diseases

Termination rate ↗

- ▶ ↗ Absenteism (diseases)

Percentage of managers ↗

- ▶ ↗ Access to training
- ▶ ↘ Shift or night working hours

Age ↗

- ▶ ↘ Fixed-term contract rate
- ▶ ↘ Productivity (weak impact)

?

- ▶ Productivity ↗ → Participation rate ↗

Global relations between QLW and performance ?

Failure

- ▶ Nothing conclusive

Interpretation

- ▶ Exist confounders (controlling QLW and performance) $C \rightarrow A$ and $C \rightarrow B$
- ▶ One such confounder is the activity sector
- ▶ In different activity sectors, causal relations are different (hampering their identification)
- ▶ \Rightarrow Condition on confounders

Low-tech sector

- ▶ Resignation rate ↗, Productivity ↘
- ▶ Average salary ↗, Productivity ↗
- ▶ Occupational-related diseases ↗, Productivity ↘
- ▶ Temporary job rate ↗, Gravity of work injuries ↗
- ▶ Permanent contract rate ↗, Safety training ↘
- ▶ Duration training stints ↗, Termination rate ↘

very significant

Outcomes & Limitations

Causal modeling and exploratory analysis

- ▶ Efficient filtering of plausible relations (several orders of magnitude);
- ▶ Complementary w.r.t. visual inspection (experts can be fooled and make sense of correlations & hazards);
- ▶ Multi-factorial relations ? yes; but even harder to interpret.

Not a ready-made analysis

- ▶ Causal relations must be
 - ▶ interpreted
 - ▶ confirmed by field experiments; polls; interviews.

Motivation

State of the art

Formal Background

The cause-effect pair challenge

The general setting

Causal Generative Neural Nets

Applications

Human Resources

Food and Health

Discussion

A data-driven approach to individual dietary recommendations

Context

- ▶ Long-term goal: Personalized dietary recommendations
- ▶ Requirement: identify risk index associated to food products
- ▶ At a coarse-grained level (lipid, protein, glucid), nothing to see
- ▶ At a fine-grained level: 300+ types of pizzas, ranging from ok to very bad.

The wealth of Kantar data

- ▶ ~22,000 households \times 10 years (this study: 2014)
- ▶ 19M total purchases/year (180,000 products)
- ▶ Socio-demographic attributes, varying size

Beware: data rarely collected as should be...

Raw description can hardly be used for meaningful analysis

- ▶ 170,000 products for 22,000 households
- ▶ Data gathered with (among others) marketing goals where bought, which conditioning
- ▶ Most products are sold by 1 vendor
- ▶ Most families are going to one vendor

Manual pre-processing

- ▶ Consider 10 categories of interest, e.g. bio/non-bio; alcohol yes/no; fresh/frozen
- ▶ Merge products with same categories
- ▶ 170,000 \rightarrow \approx 4,000 products

Example: for beer, we only selected as features of interest: colour (blonde, black, etc.); has-alcohol (yes, no); organic (yes, no)

Methodology

Dimensionality reduction

1. Borrowing Natural Language Processing tools, with
vector of purchase \approx document
food product \approx word
2. Using Latent Dirichlet Association to extract “dietary topics”

Blei et al. 03

Some topics can be directly interpreted The darker the region, the more present the topic (NB: regions are not used to build topics)



Topic 2

“Brittany”



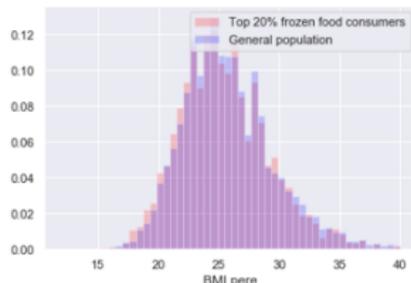
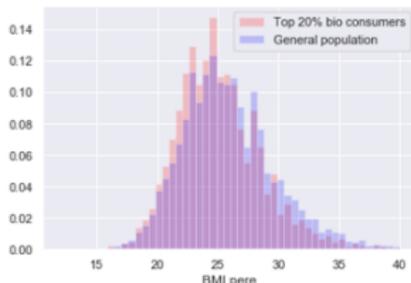
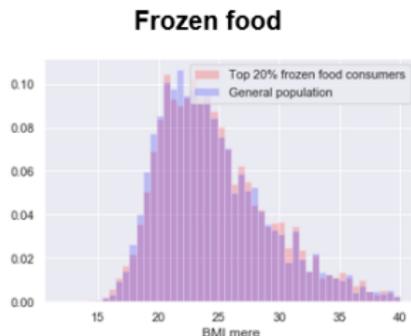
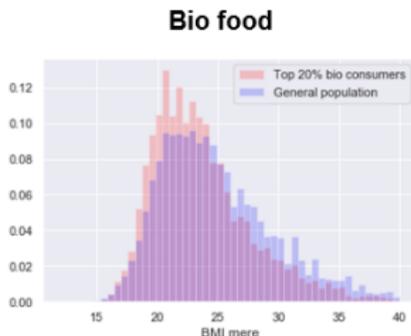
Topic 16

“Sausages++”

Focus: impact of topics on BMI

Left: Bio/organic topic
Top row: Women

Right: Frozen food topic
Bottom row: Men



High weight of Bio topic is correlated with lower BMI ($p < 5\%$)
(particularly so for women).

Does A (eat bio) cause B (better BMI) ?

Three cases

- ▶ A does cause B (bio food is better)
- ▶ Confounder: exists C that causes A and B (rich/young/educated people tend to consume bio products and have lower BMI);
- ▶ Backdoor effects: exists C correlated with A which causes B (people eating bio also tend to eat more greens, which causes lower BMI);

Goal: Find out which case holds

Causal models

- ▶ Ideally based on randomized controlled trials

Imbens Rubins 15

Proposed Methodology

Target population: "Bio" people Taking inspiration from Abadie Imbens 06
= top quantile coordinate on bio topic.

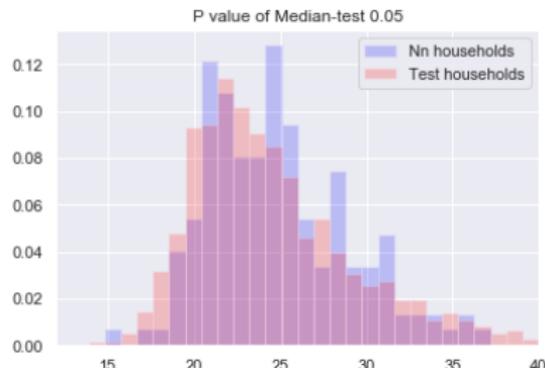
RCT would require a control population

Building a control population finding matches

- ▶ For each bio person, take her consumption z (basket of products)
- ▶ Create a falsified consumption z' (replacing each bio product with same, but non-bio, product)
- ▶ Find true consumption z'' nearest to z' (in LDA space)
- ▶ Let the true person with consumption z'' be called "falsified bio"

Compare bio and "falsified bio" populations wrt BMI

Bio vs Falsified Bio populations



Left

- ▶ Projection on the Bio topic (in log scale)
- ▶ (Falsified bio population not 0: the bio topic contains e.g. sheep yogurt).

Right

- ▶ BMI Histograms of both bio and falsified bio populations
- ▶ Statistically significant difference

Next

Chasing confounders

- ▶ Discriminating bio from “falsified bio” populations w.r.t. socio-professional features: accuracy $\approx 60\%$
- ▶ Candidate confounder: mother education level (on-going study)

Next steps

- ▶ Confirm conjectures using longitudinal data (2015-2016)
- ▶ Interact with nutritionists / sociologists
- ▶ Extend the study to consider the impact of, e.g.
 - ▶ Price of the food
 - ▶ Amount of trans fats
 - ▶ Amount of added sugar

Motivation

State of the art

Formal Background

The cause-effect pair challenge

The general setting

Causal Generative Neural Nets

Applications

Human Resources

Food and Health

Discussion

Perspectives: Causality analysis and Big Data

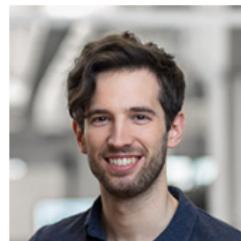
Finding the needle in the haystack

- ▶ Redundant variables (e.g. in economics) → un-interesting relations
- ▶ Variable selection
- ▶ Feature construction dimensionality reduction

Beyond causal sufficiency

- ▶ Confounders are all over the place (and many are plausible, e.g. age and size of firm; company ownership and shareholdings)
- ▶ When prior knowledge available, condition on confounders
- ▶ Use causal relationships on latent variables Wang and Blei, 19
to filter causal relationships on initial variables

Thanks!



Isabelle Guyon, Diviyam Kalainathan, Olivier Goudet, David Lopez-Paz,
Philippe Caillou, Paola Tubaro,
Ksenia Gasnikova