#### Causality, Statistics and Machine Learning

Fabio Massimo Zennaro fabio.zennaro@uib.no

University of Bergen

Bergen Meetup January 29th, 2025





#### 3 Causal Problems



Theoretically:

- It is the foundation of our understanding of the world.
- It is at the core of scientific endeavours.

Theoretically:

- It is the foundation of our understanding of the world.
- It is at the core of scientific endeavours.

Practically:

- It allows us to differentiate association and causation.
- It allows us to model non-static settings.
- It allows us to learn robust models.
- It allows us to define interventions and policies.

Theoretically:

- It is the foundation of our understanding of the world.
- It is at the core of scientific endeavours.

Practically:

- It allows us to differentiate association and causation.
- It allows us to model non-static settings.
- It allows us to learn robust models.
- It allows us to define interventions and policies.

We will follow an **operational** approach.

# 1. A Motivating Example

# Ice Creams and Thefts [9]

Assume we monitored the number of *ice-creams sold* (Ice) and the number of *thefts* (Thf) in our town:

Ice	Thf
36	20
35	18
101	31
17	12
50	23
65	25

# Ice Creams and Thefts [9]

Assume we monitored the number of *ice-creams sold* (Ice) and the number of *thefts* (Thf) in our town:

Ice	Thf
36	20
35	18
101	31
17	12
50	23
65	25

What can we infer from this data?

# The Ideal Statistician

 $\checkmark$  We learn the *joint distribution* of the variables: P(Ice, Thf)

# The Ideal Statistician

✓ We learn the *joint distribution* of the variables: P(*lce*, Thf)
✓ We can *marginalize* and *condition*: P(Thf), P(Thf|*lce*)

# The Ideal Statistician

✓ We learn the *joint distribution* of the variables: P(*lce*, *Thf*)
✓ We can *marginalize* and *condition*: P(*Thf*), P(*Thf*|*lce*)





#### The Ideal Machine Learner

 $\checkmark\,$  We can learn how the variables are correlated: Ice  $\uparrow,$  Thf  $\uparrow\,$ 

# The Ideal Machine Learner

- $\checkmark\,$  We can learn how the variables are correlated: Ice  $\uparrow,$  Thf  $\uparrow\,$
- ✓ We can *predict* a variable from another: Thf = f(Ice), Ice = f(Thf)

# The Ideal Machine Learner

- $\checkmark\,$  We can learn how the variables are correlated: Ice  $\uparrow,$  Thf  $\uparrow\,$
- ✓ We can *predict* a variable from another: Thf = f(Ice), Ice = f(Thf)





# Let's Intervene!

We might now look at these models, and try to take advantage of them:



# Let's Intervene!

We might now look at these models, and try to take advantage of them:



So, what if we stop the sale of ice-creams?

# The Naive Statistician

Let's compute the conditional for lce = 0.



# The Naive Statistician

Let's compute the conditional for lce = 0.



# The Naive Machine Learner

Let's use our model to compute lce = 0.



# The Naive Machine Learner

Let's use our model to compute lce = 0.



$$--- Thf = 3 * \sqrt{lce} + 1$$

$$Thf = 3 * \sqrt{0} + 1$$
  
 $Thf = 1$ 

# Let's Collect Data!

Let us check our conclusions against reality.

# Clashing with Reality

#### The naive answers:



# Clashing with Reality

The naive answers:



#### Collected data:

Ice	Thf
0	6
0	29
0	9
0	10
0	17
0	12
0	14

A Motivating Example

# Clashing with Reality

0.35





Collected data:



E[Thf] = 17.628

# What's the Problem in What We Did?

From the point of view of the *data model*:



- Changing *Ice* means changing the joint distribution.
- Samples are not from the same distribution anymore.

# What's the Problem in What We Did?

From the point of view of the *learned model*:

$$---- Thf = 3 * \sqrt{lce} + 1$$

- The input-output relation is not causal.
- We learned to predict a correlation, not a causal mechanism.

# Statistics/ML vs Causality [10, 6]

There are ideas we can not express in statistical/ML language.

Statistics/ML	Causality
Association	Cause
Correlation	Causation
Non-directionality	Directionality
Prediction	Action
Observation	Intervention

There is a chasm between statistics and causality.

# Questions We Can Not Express

There are questions we can not express in statistical/ML language!

Causality	3. Counterfactuals	What would have Y been, had X been x' when instead it was x? $P(Y_{do(X=x')} Y = y, X = x)$ Structural causal models
	2. Causal Effects	What is the effect of X on Y? P(Y do(X = x))
		Causal Bayesian networks
at/ML	1. Associative Relationships	How does Y relate to X?
		P(Y X)
Sta		Bayesian networks

This constitutes the Pearl's Causality Ladder [11, 12, 19, 14]

# 2. Structural Causal Models

# How to Account for Intervening?

✓ We want to learn a causal mechanism:

Effect = f(Cause)

P (Effect|Cause)

# How to Account for Intervening?

 $\checkmark$  We want to learn a causal mechanism:

Effect = f(Cause)

P (Effect|Cause)

✓ We need an idea of *directionality* between variables:



# How to Account for Intervening?

✓ We want to learn a causal mechanism:

Effect = f(Cause)

P (Effect|Cause)

✓ We need an idea of *directionality* between variables:



✓ We need to understand how correlated variables can be causally related.

#### Reichenbach's Principle

Two correlated variables X and Y can be causally related in only three ways<sup>1</sup>:  $X \rightarrow Y$ ,  $X \leftarrow Y$ ,  $X \leftarrow Z \rightarrow Y$ .

<sup>&</sup>lt;sup>1</sup>Excluding colliders and coincidences.

# Reichenbach's Principle

Two correlated variables X and Y can be causally related in only three ways<sup>1</sup>:  $X \rightarrow Y$ ,  $X \leftarrow Y$ ,  $X \leftarrow Z \rightarrow Y$ .

There likely is a *common cause* (Z) between the variables, such as the temperature:



We have a **confounder** between *Ice* and *Thf*.

<sup>&</sup>lt;sup>1</sup>Excluding colliders and coincidences.

# SCMs

**Structural causal models** provide a way to deal with interventions and counterfactuals.



We have a probabilistic model expressed via a reparametrization trick.
An intervention is a new operation do(X = x) by which a variable is set to a fixed value.

An intervention is a new operation do(X = x) by which a variable is set to a fixed value.



An intervention is a new operation do(X = x) by which a variable is set to a fixed value.



An intervention is a new operation do(X = x) by which a variable is set to a fixed value.



We obtained the new intervened (or post-intervention) model.

### Back to Our Example

We learned in an *observational* environment:



### Back to Our Example

We learned in an *observational* environment:

We deployed in this *interventional* environment:



Structural Causal Models

## (Behind the Scene: The Actual SCM in Our Example)



Structural Causal Models

# Interventions are not Conditioning

**Conditioning**  $\neq$  **Intervention** 

Structural Causal Models

### Interventions are not Conditioning

**Conditioning**  $\neq$  **Intervention** 



P(Thf|Ice = 0)

Distribution of Thf when observing Ice = 0.

Knowledge of Ice = 0 allows inference on distribution of Z and then Thf.

### Interventions are not Conditioning

**Conditioning**  $\neq$  **Intervention** 





$$P(Thf | Ice = 0)$$

Distribution of Thf when observing Ice = 0.

Knowledge of Ice = 0 allows inference on distribution of Z and then Thf.

$$P(Thf|do(X=0))$$

Distribution of Thf when intervening to do Ice = 0.

Knowledge of do (Ice = 0) does not affect the distribution of Z.

## 3. Causal Problems

### Causal Inference

Most of our data are statistical/observational data:



### Causal Inference

Most of our data are statistical/observational data:

 $\begin{array}{c} {\rm causality}\\ {\rm formalism} &\longrightarrow \begin{array}{c} {\rm statistical}\\ {\rm formalism} \end{array}$   $\begin{array}{c} {\rm interventional}\\ {\rm domain} &\longrightarrow \begin{array}{c} {\rm observational}\\ {\rm domain} \end{array}$ 

*Causal inference* provides theory and methods to exploit graphs and data to reduce *interventional queries* to *observational queries*.

#### Intervention ~> Conditioning



Given observational data, can we identify the graphical causal model  ${\cal M}$  that generated the data?

# Graph discovery [14, 8]

Given observational data, can we identify the graphical causal model  ${\cal M}$  that generated the data?

- For each probabilistic SCM there is a *single* pdf underlying it.
- For each pdf there is a set of SCMs encoding it (Markov equivalence class)

# Graph discovery [14, 8]

Given observational data, can we identify the graphical causal model  ${\cal M}$  that generated the data?

- For each probabilistic SCM there is a *single* pdf underlying it.
- For each pdf there is a set of SCMs encoding it (Markov equivalence class)

*Causal discovery* studies how to exploit data to reconstruct *casual structures*.

#### Other Causal Problems

- Learning with hidden confounders
- Causal modelling in *time-varying settings*
- Mediation analysis
- Inference with missing data
- Inference with partially specified models
- Discovery with interventions
- Experimental design
- Causal transportability
- Counterfactual reasoning

• ...

#### Relation to Machine Learning

A double relation: ML can use causality theory to improve learning, and causaliy theory can use ML to improve causal inference.

Sample intersections:

- Causal and anti-causal learning
- Invariance learning
- Reinforcement learning
- Counterfactual fairness
- Causal abstraction learning

### Advantages of causality

Causal reasoning is not necessary if:

• We want to model/predict in a static setting.

Causal reasoning is not necessary if:

• We want to model/predict in a static setting.

However, causal modelling may allow us (among other things) to:

- Distinguish and learn actual causal mechanisms;
- Deal with settings changing under *interventions*.

Causal reasoning is not necessary if:

• We want to model/predict in a static setting.

However, causal modelling may allow us (among other things) to:

- Distinguish and learn actual causal mechanisms;
- Deal with settings changing under *interventions*.

(Causal libraries are available, such as *do-why* or *causal-learn*)

The theory of causality empowers machine learning:

- Provides a formalism to reason causally (the SCM framework is general, it helps making assumptions explicit, and it eases reasoning via graphs).
- Allows to express causal statements.
- Allows for learning *robust* models.
- Enhance *interpretability* and *explainability* of models.
- It may spur us to move beyond deep learning.

It comes with a **cost** though:

• Assumptions/structures!

### Thanks!

Thank you for listening!

#### References I

- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. arXiv preprint arXiv:1907.02893, 2019.
- [2] Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in criminal justice risk assessments: the state of the art. arXiv preprint arXiv:1703.09207, 2017.
- [3] Léon Bottou, Jonas Peters, Joaquin Quiñonero-Candela, Denis X Charles, D Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning systems: The example of computational advertising. *The Journal of Machine Learning Research*, 14(1):3207–3260, 2013.
- [4] Lars Buesing, Theophane Weber, Yori Zwols, Sebastien Racaniere, Arthur Guez, Jean-Baptiste Lespiau, and Nicolas Heess. Woulda, coulda, shoulda: Counterfactually-guided policy search. arXiv preprint arXiv:1811.06272, 2018.

#### References II

- [5] Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163, 2017.
- [6] A Philip Dawid. Statistical causality from a decision-theoretic perspective. Annual Review of Statistics and Its Application, 2:273–303, 2015.
- [7] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In Advances in Neural Information Processing Systems, pages 4069–4079, 2017.
- [8] Marloes H Maathuis, Preetam Nandy, and P Btihlmann. A review of some recent advances in causal inference., 2016.
- [9] Judea Pearl. Causality. Cambridge University Press, 2009.
- [10] Judea Pearl. An introduction to causal inference. *The international journal of biostatistics*, 2010.

#### References III

- [11] Judea Pearl. Theoretical impediments to machine learning with seven sparks from the causal revolution. arXiv preprint arXiv:1801.04016, 2018.
- [12] Judea Pearl. Sufficient causes: Revisiting oxygen, matches, and fires. *Journal of Causal Inference*, 2019.
- [13] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 78(5):947–1012, 2016.
- [14] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Elements of causal inference: Foundations and learning algorithms. MIT Press, 2017.
- [15] Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for causal transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342, 2018.

#### References IV

- [16] Chris Russell, Matt J Kusner, Joshua Loftus, and Ricardo Silva. When worlds collide: integrating different counterfactual assumptions in fairness. In Advances in Neural Information Processing Systems, pages 6417–6426, 2017.
- [17] Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris Mooij. On causal and anticausal learning. arXiv preprint arXiv:1206.6471, 2012.
- [18] Bernhard Schölkopf. Causality for machine learning. *arXiv preprint arXiv:1911.10500*, 2019.
- [19] Ilya Shpitser and Judea Pearl. Complete identification methods for the causal hierarchy. *Journal of Machine Learning Research*, 9(Sep):1941–1979, 2008.
- [20] Guy Tennenholtz, Shie Mannor, and Uri Shalit. Off-policy evaluation in partially observable environments. arXiv preprint arXiv:1909.03739, 2019.

A SCM expresses and encodes statistical and causal assumptions:

- Acyclicity: no loops in the graph.
- *Causal Markov assumption*: a node is independent of its non-effects given its direct causes.
- Zero influence: missing arrow means no causal relationship.
- Common cause completeness: all common causes are modeled.
- *Autonomous functions*: changing a function does not affect other functions.

...

#### No causes in, no causes out.

### Structural causal models [18]

P(X, Y)



- Statistics works with the *joint*; factorizations are instrumental.
- Causality makes the *assumption* that one of the factorizations is the *true causal model*.

#### A causal model contains more information than a statistical one.

#### A Motivating Example

SCMs represent causal systems.



### A Motivating Example

SCMs represent causal systems.



SCMs integrates a graphical model and probabilities distributions.

## Structural Causal Models (SCMs) - Definition

We express a **SCM** as  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$  [9, 14]:



## Structural Causal Models (SCMs) - Definition

We express a **SCM** as  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$  [9, 14]:



• X: set of *endogenous nodes* (S, T, C) representing variables of interest

## Structural Causal Models (SCMs) - Definition

We express a **SCM** as  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$  [9, 14]:



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U<sub>S</sub>, U<sub>T</sub>, U<sub>C</sub>) representing stochastic factors

### Structural Causal Models (SCMs) - Definition

We express a **SCM** as  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$  [9, 14]:



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U<sub>S</sub>, U<sub>T</sub>, U<sub>C</sub>) representing stochastic factors
- *F*: Set of *structural functions* (*f<sub>S</sub>*, *f<sub>T</sub>*, *f<sub>C</sub>*) describing the dynamics of each variable
## Structural Causal Models (SCMs) - Definition

We express a **SCM** as  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$  [9, 14]:



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U<sub>S</sub>, U<sub>T</sub>, U<sub>C</sub>) representing stochastic factors
- *F*: Set of *structural functions* (*f<sub>S</sub>*, *f<sub>T</sub>*, *f<sub>C</sub>*) describing the dynamics of each variable
- $\mathcal{P}$ : Set of *distributions* ( $P_S, P_T, P_C$ ) describing the random factors

## Structural Causal Models (SCMs) - Definition

We express a **SCM** as  $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$  [9, 14]:



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U<sub>S</sub>, U<sub>T</sub>, U<sub>C</sub>) representing stochastic factors
- *F*: Set of *structural functions* (*f<sub>S</sub>*, *f<sub>T</sub>*, *f<sub>C</sub>*) describing the dynamics of each variable
- $\mathcal{P}$ : Set of *distributions* ( $P_S, P_T, P_C$ ) describing the random factors

Every SCM  $\mathcal{M}$  implies a (joint) distribution  $P_{\mathcal{M}}$ :  $P_{\mathcal{M}}(S, T, C)$ 

F.M. Zennaro

# Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [9, 14]:





# Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [9, 14]:



do(T = 1)

# Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [9, 14]:



do(T=1)

2

Remove incoming edges in the intervened node

# Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [9, 14]:



do(T=1)

- Remove incoming edges in the intervened node
- Set the value of the intervened node





An *intervention*  $\iota$  defines a new **intervened model**  $\mathcal{M}_{\iota}$  with new distributions.



 $P_{\mathcal{M}}$ 







A **counterfactual** is an operation by which we compute a quantity of interest in an alternate world in which we perform an intervention.

$$P\left(Y_{do(X=x')}|Y=y,X=x\right)$$

This reflects the *counterfactual question*: assuming we observed Y = y and X = x, what would have Y been, had we acted on do(X = x')?

Interventions  $\neq$  Counterfactuals



$$P(Bet = Coin | do(Bet = head))$$

Probability of winning if we force the bet to head.

The outcome of the coin toss is still random, and the chance of winning half.

$$P(Bet = Coin_{do(Bet = head)}|$$
  
 $Coin = head, Bet = tail)$ 

Probability of winning if we had forced the bet to head, having observed the outcome head and the bet tail.

We know with certainty the result of the bet.

## Causal and Anti-Causal Learning [17]

# Causal Learning

Given samples (cause, effect) we learn:

Effect = f(Cause)

P (Effect|Cause)

**Anti-Causal Learning** 

Given samples *(effect, cause)* we learn:

 $\mathsf{Cause} = f(\mathsf{Effect})$ 

P(Cause|Effect)

e.g.: predicting structure of proteins.

e.g.: classifying images.

 $P(\mathsf{Effect}|\mathsf{Cause}) \perp P(\mathsf{Cause})$ 

# Semi-supervised Learning [17]

## $P(\mathsf{Effect}|\mathsf{Cause}) \perp P(\mathsf{Cause})$

### **Causal Learning**

In SSL, we receive more samples *(cause)*, and we aim to learn:

P (Effect|Cause)

Learning more on how the cause distributes do not provide information on how the effect mechanism behaves. (But it may help reducing the risk!)

### **Anti-Causal Learning**

In SSL, we receive more samples *(effect)*, and we aim to learn:

### P (Cause|Effect)

Learning more on how the effect distributes may help us infer more about the cause mechanism under standard SSL assumptions (smoothness, clustering).

## Covariate Shift [17]

### $P(\mathsf{Effect}|\mathsf{Cause}) \perp P(\mathsf{Cause})$

### **Causal Learning**

In CS, we receive test samples from P'(Cause), and we aim to compute:

P (Effect|Cause)

The effect mechanism is not affected by shifts in the distribution of the causes. (But risk may require adjustment!)

#### Anti-Causal Learning

In CS, we receive test samples from P'(Effect), and we aim to compute:

P(Cause|Effect)

A change in the effect mechanism affects the conditional distribution of causes.

In absence of a model, we may try to learn a *local structure*.

Suppose we are given data from different *environment* (= *interventional domains*)



## Invariance Learning

From data in different environments  $\mathcal{E}_i$  we can learn the sets of variables that is *invariant* in all the settings (= under all interventions).



The set of invariant variables are the (true) *direct causes* of the variable of interest.

Prediction of invariance [13, 15] and learning with invariant risk minimization [1] allow for learning robust model (= transfer learning).

Reinforcement learning deals with an interventional setting.

Performing actions, an agent probes the distribution of an environment under intervention:

$$P(E|do(A = a))$$

*Bandit problems* and *reinforcement learning* may be expressed in causal terms [3].

Reinforcement learning works without structural models and causal formalization is still debated.

There are promising point of contacts:

- Counterfactual reasoning with structure in ad placement problems [3]
- Relation between *offline policy evaluation* and *inverse probability weighting* [3, 20]
- Counterfactually-guided policy search [4]



Fairness is concerned with deciding if learned systems are socially fair.



Important in applications such as job recruiting, loan decisions, police deployment.

## How to Measure Fairness?

There are several approaches to guarantee fairness [2]:

- Fairness through unawareness:  $\hat{Y} = f(X)$
- Demographic parity:  $P(\hat{Y}|A=0) = P(\hat{Y}|A=1)$
- Equality of opportunity:  $P(\hat{Y}|A=0, Y=1) = P(\hat{Y}|A=1, Y=1)$

These measures are either insufficient [7] or conflicting [5].

# Counterfactual Fairness [7, 16]

We can enforce an individual-level fairness in *counterfactual* terms:

$$P\left(\hat{Y}|X=x,A=a\right) = P\left(\hat{Y}_{do(A=a')}|X=x,A=a\right)$$

For instance:

$$\begin{array}{l} P(\text{accepted}|X = x, A = \text{female}) \\ = \\ P\left(\text{accepted}_{do(A = \text{male})} | X = x, A = \text{female}\right) \end{array}$$