Causal Inference from Observational Data Based on Graphical Models

Guido Consonni

O'Bayes 2025, June 8, Athens

Outline



From Statistical to Causal Models

2 Causal Modeling

3 Causal Discovery





Turing award – 2011 Judea Pearl

Turing award – 2011 Judea Pearl "For fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning"



Turing award – 2011

Judea Pearl

"For fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning" The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2021 (David Card) Joshua D. Angrist and Guido W.

Imbens



Turing award – 2011

Judea Pearl

"For fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning"



The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2021 (David Card) Joshua D. Angrist and Guido W. Imbens "For their methodological contributions to the analysis of causal relationships"



Guido Consonni

O'Bayes 2025, June 8, Athens

Rousseeuw Prize for Statistics 2022

James Robins, Miguel Hernán, Thomas Richardson, Andrea Rotnitzky, Eric Tchetgen Tchetgen

"For their pioneering work on Causal Inference with applications in Medicine and Public Health"



"The kind of causal inference seen in natural human thought can be "algorithmitized" to help produce human-level machine intelligence" Judea Pearl, 2019, *Communications of the ACM* "The kind of causal inference seen in natural human thought can be "algorithmitized" to help produce human-level machine intelligence" Judea Pearl, 2019, *Communications of the ACM*

"Some of the hard open problems of machine learning and AI are intrinsically related to causality, and progress may require advances in our understanding of how to model and infer causality from data" Bernhard Schölkopf, 2022, International Congress of Mathematicians "The kind of causal inference seen in natural human thought can be "algorithmitized" to help produce human-level machine intelligence" Judea Pearl, 2019, *Communications of the ACM*

"Some of the hard open problems of machine learning and AI are intrinsically related to causality, and progress may require advances in our understanding of how to model and infer causality from data" Bernhard Schölkopf, 2022, International Congress of Mathematicians

Preserving Causal Constraints in Counterfactual Explanations for Machine Learning Classifiers *NeurIPS*, 2019

Interpretable ML; feasible counterfactuals

Preserving Causal Constraints in Counterfactual Explanations for Machine Learning Classifiers *NeurIPS*, 2019 Interpretable ML; feasible counterfactuals

Causal Reasoning and Large Language Models: Opening a New Frontier for Causality

Transactions on Machine Learning Research, 2024

"Behavorial" study of LLMs to benchmark their capability in generating causal arguments

Preserving Causal Constraints in Counterfactual Explanations for Machine Learning Classifiers *NeurIPS*, 2019 Interpretable ML; feasible counterfactuals

Causal Reasoning and Large Language Models: Opening a New Frontier for Causality

Transactions on Machine Learning Research, 2024

"Behavorial" study of LLMs to benchmark their capability in generating causal arguments

Improving the accuracy of medical diagnosis with causal machine learning *Nature communications*, 2020

we reformulate diagnosis as a counterfactual inference task and derive counterfactual diagnostic algorithms.In medical diagnosis a doctor aims to explain a patient's symptoms by determining the diseases causing them, while existing diagnostic algorithms are purely associative

Robust Agents Learn Causal World Models International Conference on Learning Representations, 2024

"Any agent capable of satisfying a regret bound for a large set of distributional shifts must have learned an approximate causal model of the data generating process"

Robust Agents Learn Causal World Models International Conference on Learning Representations, 2024

"Any agent capable of satisfying a regret bound for a large set of distributional shifts must have learned an approximate causal model of the data generating process"

Explaining the Behavior of Black-Box Prediction Algorithms with Causal Learning

Transactions on Machine Learning Research, 2025

Causal approaches to post-hoc explainability for black-box prediction models(e.g. deep neural networks trained on image pixel data)

Table of Contents



From Statistical to Causal Models



Correlation does not imply causation Chocolate and Nobel prize winners



Correlation does not imply causation Chocolate and Nobel prize winners



Understanding causation

- Manipulability
- Intervention

Correlation does not imply causation Chocolate and Nobel prize winners



Understanding causation

- Manipulability
- Intervention

J. Woodward (2001). *Causation and manipulability* J. Pearl (2009). *Causality: models, reasoning, and inference*. 2nd edn

Correlation does not imply causation Chocolate and Nobel prize winners



- Understanding causation
 - Manipulability
 - Intervention

J. Woodward (2001). *Causation and manipulability* J. Pearl (2009). *Causality: models, reasoning, and inference*. 2nd edn

Epidemiology
 I. M. Babina M

J. M. Robins, M. A. Hernan, and B. Brumback (2000)

- Agriculture
 S. Wright (1921)
- Econometrics

T. Haavelmo (1944); K. D. Hoover (2001)

Causal Reasoning

Definition

A random variable X has a causal effect on a random variable Y if there exist $x \neq x'$ such that the distribution of Y after intervening on X and setting it to x differs from the distribution of Y after setting X to x'

Definition

A random variable X has a causal effect on a random variable Y if there exist $x \neq x'$ such that the distribution of Y after intervening on X and setting it to x differs from the distribution of Y after setting X to x'



Definition

A random variable X has a causal effect on a random variable Y if there exist $x \neq x'$ such that the distribution of Y after intervening on X and setting it to x differs from the distribution of Y after setting X to x'



Definition

A random variable X has a causal effect on a random variable Y if there exist $x \neq x'$ such that the distribution of Y after intervening on X and setting it to x differs from the distribution of Y after setting X to x'



Gene A is correlated with the phenotype, and so is gene B However only gene A ha a causal effect on the phenotype

Causal Reasoning

Correlation and Causation: what's the connection?

Principle

If two random variables X and Y are statistically dependent $X \not\perp Y$ then there exists a random variable Z which causally influences both of them and which explains all their dependence that is $X \perp Y \mid Z$ (c) As a special case, Z may coincide with X or Y (a) or (b)



(b) X (c) 🗸

(a) X



X: Chocolate consumptionY: # Nobel laureatesZ: Economic factor

(b) X (c) 🗸

(a) X



X: Chocolate consumptionY: # Nobel laureatesZ: Economic factor







• The class of observational distributions over X and Y that can be realized by these models is the same in all three cases





X: Chocolate consumptionY: # Nobel laureatesZ: Economic factor

Causal Reasoning

O'Bayes 2025, June 8, Athens

(a) \mathbf{X} (X) \rightarrow (Y)







X: Chocolate consumptionY: # Nobel laureatesZ: Economic factor

- The class of observational distributions over X and Y that can be realized by these models is the same in all three cases
- Cannot distinguish among a), b) and c) through passive observation i.e., in a purely data-driven way

(a) \mathbf{X} (X) \rightarrow (Y)





X: Chocolate consumption Y: # Nobel laureates Z: Economic factor

- The class of observational distributions over X and Y that can be realized by these models is the same in all three cases
- Cannot distinguish among a), b) and c) through passive observation i.e., in a purely data-driven way
- Z latent confounder drives consumer spending and investment in education and research [from background knowledge]

Causal Reasoning

• Correlation is still useful

- Correlation is still useful
- Causality is not always needed

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting
 - we measure the activities of certain genes and are asked to predict the phenotype
- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting
 - we measure the activities of certain genes and are asked to predict the phenotype

• However, if we want to answer interventional questions

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting
 - we measure the activities of certain genes and are asked to predict the phenotype

- However, if we want to answer interventional questions
 - the outcome of a gene knockout experiment

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting
 - we measure the activities of certain genes and are asked to predict the phenotype

- However, if we want to answer interventional questions
 - the outcome of a gene knockout experiment
 - the effect of a policy enforcing a job training program (or higher chocolate consumption)

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting
 - we measure the activities of certain genes and are asked to predict the phenotype

- However, if we want to answer interventional questions
 - the outcome of a gene knockout experiment
 - the effect of a policy enforcing a job training program (or higher chocolate consumption)
- We need more than correlation

- Correlation is still useful
- Causality is not always needed
- Gene A and gene B remain useful features for making predictions
- In a passive, or observational, setting
 - we measure the activities of certain genes and are asked to predict the phenotype

- However, if we want to answer interventional questions
 - the outcome of a gene knockout experiment
 - the effect of a policy enforcing a job training program (or higher chocolate consumption)
- We need more than correlation
- We need a causal model

Table of Contents

From Statistical to Causal Models

2 Causal Modeling

3 Causal Discovery

4 Causal Reasoning

5 Conclusions

Causal Graphical Model

Definition

A Causal Graphical Model (CGM) $\mathcal{M} = (G, p)$ over n random variables X_1, \ldots, X_n consists of

- a directed acyclic graph (DAG) G in which directed edges $(X_j \rightarrow X_i)$ represent a direct causal effect of X_j on X_i ;
- a joint distribution $p(X_1,\ldots,X_n)$ which is Markovian w.r.t. G

$$p(X_1, \dots, X_n) = \prod_{i=1} p(X_i | PA_i); \quad PA_i = \{X_j : (X_j \to X_i) \in G\}$$

Causal Graphical Model

Definition

A Causal Graphical Model (CGM) $\mathcal{M} = (G, p)$ over n random variables X_1, \ldots, X_n consists of

- a directed acyclic graph (DAG) G in which directed edges $(X_j \rightarrow X_i)$ represent a direct causal effect of X_j on X_i ;
- a joint distribution $p(X_1, \ldots, X_n)$ which is Markovian w.r.t. G $p(X_1, \ldots, X_n) = \prod_{i=1}^n p(X_i | PA_i); \quad PA_i = \{X_j : (X_j \to X_i) \in G\}$

 PA_i is the set of parents, or direct causes, of X_i in G

Causal Graphical Model

Definition

A Causal Graphical Model (CGM) $\mathcal{M} = (G, p)$ over n random variables X_1, \ldots, X_n consists of

- a directed acyclic graph (DAG) G in which directed edges $(X_j \rightarrow X_i)$ represent a direct causal effect of X_j on X_i ;
- a joint distribution $p(X_1, \ldots, X_n)$ which is Markovian w.r.t. G $p(X_1, \ldots, X_n) = \prod_{i=1}^n p(X_i | PA_i); \quad PA_i = \{X_j : (X_j \to X_i) \in G\}$

 PA_i is the set of parents, or direct causes, of X_i in GDecomposition of the joint distribution into causal conditionals

Four variables



Four variables



 $P(X_1, X_2, X_3, X_4) =$ P(X_1) P(X_4) P(X_2 | X_1) P(X_3 | X_1, X_2, X_4)

Definition

A joint distribution satisfies the causal Markov condition w.r.t. a DAG G if every variable is conditionally independent of its non-descendants in G given its parents in G

Definition

A joint distribution satisfies the causal Markov condition w.r.t. a DAG G if every variable is conditionally independent of its non-descendants in G given its parents in G



Definition

A joint distribution satisfies the causal Markov condition w.r.t. a DAG G if every variable is conditionally independent of its non-descendants in G given its parents in G



 $X_2 \perp \!\!\!\perp X_4 \mid X_1$ $X_4 \perp \!\!\!\perp \{X_1, X_2\}$

Definition

A joint distribution satisfies the causal Markov condition w.r.t. a DAG G if every variable is conditionally independent of its non-descendants in G given its parents in G



 $p(X_1,\ldots,X_n) = \prod_{i=1}^n p(X_i \,|\, PA_i)$ iff the Causal Markov Condition holds

Intervention on a causal DAG

Central idea

Intervening on a variable, by externally forcing it to take on a particular value, renders it independent of its causes

Intervention on a causal DAG

Central idea

Intervening on a variable, by externally forcing it to take on a particular value, renders it independent of its causes

and breaks their causal influence

Intervention on a causal DAG

Central idea

Intervening on a variable, by externally forcing it to take on a particular value, renders it independent of its causes

and breaks their causal influence

- do-operator
- graph-surgery

Three variables and a graph



Three variables and a graph



From graph G to G'



Starting graph G

Guido Consonni

From graph G to G'



Post-intervention graph G' for $do(X_2 = x_2)$.

Guido Consonni

From graph G to G''



Starting graph ${\cal G}$

Guido Consonni

From graph G to G''



Post-intervention graph G'' for $do(X_3 = x_3)$.

• Intervention and Conditioning radically different

- Intervention and Conditioning radically different
- Conditioning is passive

- Intervention and Conditioning radically different
- Conditioning is passive
- Intervention is active

- Intervention and Conditioning radically different
- Conditioning is passive
- Intervention is active
 - if a gene is knocked out, it is no longer influenced by other genes that were previously regulating it

- Intervention and Conditioning radically different
- Conditioning is passive
- Intervention is active
 - if a gene is knocked out, it is no longer influenced by other genes that were previously regulating it

instead, its activity is now solely determined by the intervention

- Intervention and Conditioning radically different
- Conditioning is passive
- Intervention is active
 - if a gene is knocked out, it is no longer influenced by other genes that were previously regulating it

instead, its activity is now solely determined by the intervention

Note

This is fundamentally different from conditioning, because passively observing the activity of a gene provides information about its driving factors (i.e., its direct causes)

- Intervention and Conditioning radically different
- Conditioning is passive
- Intervention is active
 - if a gene is knocked out, it is no longer influenced by other genes that were previously regulating it

instead, its activity is now solely determined by the intervention

Note

This is fundamentally different from conditioning, because passively observing the activity of a gene provides information about its driving factors (i.e., its direct causes)

$$p(y \mid x) \neq p(y \mid do(X = x))$$





Example



$$p(X_3|do(X_2 = x_2)) = \sum_{x_1} p(x_1)p(X_3|x_1, x_2)$$

Guido Consonni

Example



 $p(X_3|do(X_2 = x_2)) = \sum_{x_1} p(x_1)p(X_3|x_1, x_2)$



Example



$$p(X_3|do(X_2 = x_2)) = \sum_{x_1} p(x_1)p(X_3|x_1, x_2)$$



$$p(X_3 \mid x_2) = \sum_{x_1} p(x_1 \mid x_2) p(X_3 \mid x_1, x_2)$$

Definition

An SCM $\mathcal{M} = (F, p_U)$ consists of

i) a set F of n assignments (the structural equations)

$$F = \{X_i := f_i(PA_i, U_i), i = 1, \dots, n\}$$

 $PA_i \subseteq \{X_1, \dots, X_n\} \setminus X_i$: causal parents U_i 's: noise variables

ii) a joint distribution $p_U(U_1, \ldots, U_n)$
Features of an SCM

- Each X_i is generated from other variables through a deterministic mechanism ${\cal F}$

- Each X_i is generated from other variables through a deterministic mechanism F
- Randomness originates from U_i's stochasticity of the process uncertainty due to unmeasured parts

- Each X_i is generated from other variables through a deterministic mechanism F
- Randomness originates from U_i 's stochasticity of the process uncertainty due to unmeasured parts
- $X_i := f_i(PA_i, U_i)$ asymmetry between LHS and RHS

- Each X_i is generated from other variables through a deterministic mechanism F
- Randomness originates from U_i's stochasticity of the process uncertainty due to unmeasured parts
- $X_i := f_i(PA_i, U_i)$ asymmetry between LHS and RHS
- In parametric linear form (linear f_i)
 SCMs are also known as structural equation models (path analysis)

Linking SCM's and CGM's

Definition

The causal graph G induced by an SCM is the directed graph with vertex set $\{X_1, \ldots, X_n\}$ and a directed edge from each vertex in PA_i to X_i for all i.

Linking SCM's and CGM's

Definition

The causal graph G induced by an SCM is the directed graph with vertex set $\{X_1, \ldots, X_n\}$ and a directed edge from each vertex in PA_i to X_i for all i.

Example SCM over $\{X_1, X_2, X_3\}$ with some $p_U(U_1, U_2, U_3)$

$$X_1 := f_1(U_1), X_2 := f_2(X_1, U_2), X_3 := f_3(X_1, X_2, U_3)$$

Linking SCM's and CGM's

Definition

The causal graph G induced by an SCM is the directed graph with vertex set $\{X_1, \ldots, X_n\}$ and a directed edge from each vertex in PA_i to X_i for all i.

Example SCM over $\{X_1, X_2, X_3\}$ with some $p_U(U_1, U_2, U_3)$

$$X_1 := f_1(U_1), X_2 := f_2(X_1, U_2), X_3 := f_3(X_1, X_2, U_3)$$



Difference between ${\cal CGM}$ and ${\cal SCM}$

SCM allows for a rich class of causal models including models with cyclic causal relations not obeying the causal Markov condition (because of complex covariance structures between the noise terms)

Difference between ${\cal CGM}$ and ${\cal SCM}$

SCM allows for a rich class of causal models including models with cyclic causal relations *not* obeying the causal Markov condition (because of complex covariance structures between the noise terms)

Further common assumptions

- A1) Acyclicity: the induced graph G is a DAG
- A2) Causal sufficiency/no hidden confounders: the U_i 's are jointly independent, i.e.

$$p_U(U_1,\ldots,U_n)=p_{U_1}(U_1)\times\ldots p_{U_n}(U_n)$$

Difference between ${\cal CGM}$ and ${\cal SCM}$

SCM allows for a rich class of causal models including models with cyclic causal relations *not* obeying the causal Markov condition (because of complex covariance structures between the noise terms)

Further common assumptions

- A1) Acyclicity: the induced graph G is a DAG
- A2) Causal sufficiency/no hidden confounders: the U_i 's are jointly independent, i.e.

$$p_U(U_1,\ldots,U_n)=p_{U_1}(U_1)\times\ldots p_{U_n}(U_n)$$

Acyclicity and Causal sufficiency ensure that the distribution induced by an SCM *factorises* according to its induced causal graph G (and the causal Markov condition is satisfied w.r.t. G)

Definition

An intervention $do(X_i = x_i)$ in an $SCM \mathcal{M} = (F, p_U)$ is modeled by

- replacing the *i*-th structural equation in F by $X_i = x_i$
- remaining F_j 's remain unchanged $(j \neq i)$

Result is the interventional $SCM \mathcal{M}^{do(X_i=x_i)} = (F', p_U).$

Definition

An intervention $do(X_i = x_i)$ in an $SCM \mathcal{M} = (F, p_U)$ is modeled by

- replacing the *i*-th structural equation in F by $X_i = x_i$
- remaining F_j 's remain unchanged $(j \neq i)$

Result is the interventional $SCM \mathcal{M}^{do(X_i=x_i)} = (F', p_U).$

From $\mathcal{M}^{do(X_i=x_i)} = (F', p_U)$ deduce the interventional distribution $p(X_{-i} | do(X_i = x_i))$ and the intervention graph G'

Interventions in SCM

 $SCM \mathcal{M} = (F, p_U)$

 $X_1 := f_1(U_1), X_2 := f_2(X_1, U_2), X_3 := f_3(X_1, X_2, U_3)$

¹This way of handling interventions coincides with that for CGMs

Causal Reasoning

Interventions in SCM

 $SCM \mathcal{M} = (F, p_U)$

 $X_1 := f_1(U_1), X_2 := f_2(X_1, U_2), X_3 := f_3(X_1, X_2, U_3)$

 $SCM \ \mathcal{M}^{do(X_2=x_2)} = (F', p_U)$

$$X_1 := f_1(U_1), X_2 := x_2, X_3 := f_3(X_1, X_2, U_3)$$

¹This way of handling interventions coincides with that for CGMs

Interventions in SCM

 $SCM \mathcal{M} = (F, p_U)$

 $X_1 := f_1(U_1), X_2 := f_2(X_1, U_2), X_3 := f_3(X_1, X_2, U_3)$

 $SCM \ \mathcal{M}^{do(X_2=x_2)} = (F', p_U)$

$$X_1 := f_1(U_1), X_2 := x_2, X_3 := f_3(X_1, X_2, U_3)$$

Graph G' induced by $\mathcal{M}^{do(X_2=x_2)}$ 1



¹This way of handling interventions coincides with that for CGMs

Seeing, Doing, Imagining The ladder of causality

Seeing, Doing, Imagining The ladder of causality



Seeing, Doing, Imagining The ladder of causality





i) observation

i) observation

passively seen or measured

i) observation

passively seen or measured

ii) intervention

i) observation

passively seen or measured

ii) intervention

external manipulation or experimentation

i) observation

passively seen or measured

ii) intervention

external manipulation or experimentation

iii) counterfactual

i) observation

passively seen or measured

ii) intervention

external manipulation or experimentation

iii) counterfactual

what would have been, given that something else was in fact observed

i) observation

passively seen or measured

ii) intervention

external manipulation or experimentation

iii) counterfactual

what would have been, given that something else was in fact observed

Issues with counterfactuals

Cannot be observed empirically <u>unfalsifiable</u>

i) observation

passively seen or measured

ii) intervention

external manipulation or experimentation

iii) counterfactual

what would have been, given that something else was in fact observed

Issues with counterfactuals

Cannot be observed empirically unfalsifiable

<u>unscientific</u> (Popper, 1959) problematic (Dawid, 2000)

i) observation

passively seen or measured

ii) intervention

external manipulation or experimentation

iii) counterfactual

what would have been, given that something else was in fact observed

Issues with counterfactuals

Cannot be observed empirically <u>unfalsifiable</u>

unscientific (Popper, 1959) problematic (Dawid, 2000)

Yet, humans seem to perform counterfactual reasoning in practice starting in early childhood (Buchsbaum et al., 2012)

Guido Consonni

Causal Reasoning

"Given that patient X received treatment A and their health got worse, what would have happened if they had been given treatment B instead, *all else being equal*?"

"Given that patient X received treatment A and their health got worse, what would have happened if they had been given treatment B instead, *all else being equal*?"

· SCMs provide a suitable framework for counterfactual reasoning

"Given that patient X received treatment A and their health got worse, what would have happened if they had been given treatment B instead, *all else being equal*?"

- SCMs provide a suitable framework for counterfactual reasoning
- Observing what actually happened provides information about the *background state* of the system namely the noise terms $\{U_1, \ldots, U_n\}$ in an SCM

"Given that patient X received treatment A and their health got worse, what would have happened if they had been given treatment B instead, *all else being equal*?"

- SCMs provide a suitable framework for counterfactual reasoning
- Observing what actually happened provides information about the *background state* of the system namely the noise terms $\{U_1, \ldots, U_n\}$ in an SCM
- This differs from an intervention where such background information is not available

• Observing that treatment A did not work may tell us that the patient has a rare condition

- Observing that treatment A did not work may tell us that the patient has a rare condition
 - this provides information on their background state of health

- Observing that treatment A did not work may tell us that the patient has a rare condition
 - this provides information on their background state of health
- This, in turn, suggests that treatment *B* might have worked

- Observing that treatment A did not work may tell us that the patient has a rare condition
 - this provides information on their background state of health
- This, in turn, suggests that treatment *B* might have worked

- However, given that treatment A has been applied, patient's condition may have changed
 - so condition "all else being equal" fails
- Observing that treatment A did not work may tell us that the patient has a rare condition
 - this provides information on their background state of health
- This, in turn, suggests that treatment *B* might have worked

- However, given that treatment A has been applied, patient's condition may have changed
 - so condition "all else being equal" fails
- and *B* may no longer work in a future intervention *on this specific patient*

Definition (Counterfactuals in *SCM*'s)

Given evidence X = x observed from an $SCM \ \mathcal{M} = (F, p_U)$ the counterfactual $SCM \ \mathcal{M}^{X=x}$ is obtained by updating p_U to $p_{U|X=x}$

$$\mathcal{M}^{X=x} = (F, p_{U \mid X=x})$$

Definition (Counterfactuals in *SCM*'s)

Given evidence X = x observed from an $SCM \ \mathcal{M} = (F, p_U)$ the counterfactual $SCM \ \mathcal{M}^{X=x}$ is obtained by updating p_U to $p_{U|X=x}$

$$\mathcal{M}^{X=x} = (F, p_{U \mid X=x})$$

Counterfactuals are then computed by performing interventions in the counterfactual $SCM \ \mathcal{M}^{X=x}$

$SCM \mathcal{M} = (F, p_U)$

 $SCM \mathcal{M} = (F, p_U)$

 $X := U_X, Y := 3X + U_Y; U_X, U_Y \stackrel{iid}{\sim} N(0, 1)$

 $SCM \mathcal{M} = (F, p_U)$

$$X := U_X, Y := 3X + U_Y; U_X, U_Y \stackrel{iid}{\sim} N(0, 1)$$

We observe X = 2 and Y = 6.5and want to answer the counterfactual question "What would Y have been, had X = 1?"

 $SCM \mathcal{M} = (F, p_U)$

$$X := U_X, Y := 3X + U_Y; U_X, U_Y \stackrel{iid}{\sim} N(0, 1)$$

We observe X = 2 and Y = 6.5and want to answer the counterfactual question "What would Y have been, had X = 1?"

We are thus interested in

$$p^{\mathcal{M}^{X=2,Y=6.5;do(X=1)}}(Y)$$

Recall: $X := U_X$, $Y := 3X + U_Y$, U_X , $U_Y \stackrel{iid}{\sim} N(0,1)$

Recall: $X := U_X$, $Y := 3X + U_Y$, U_X , $U_Y \stackrel{iid}{\sim} N(0,1)$

• Update the noise distribution $p_U \rightarrow p_{U\,|\,X=2,Y=6.5}$

 $U_X \sim \delta(2), U_Y \sim \delta(0.5)$

Recall: $X := U_X$, $Y := 3X + U_Y$, U_X , $U_Y \stackrel{iid}{\sim} N(0,1)$

- Update the noise distribution $p_U \rightarrow p_{U\,|\,X=2,Y=6.5}$

 $U_X \sim \delta(2), U_Y \sim \delta(0.5)$

- Obtain the updated SCM $\mathcal{M}^{X=2,Y=6.5}=(F,p_{U\,|\,X=2,Y=6.5})$

Recall: $X := U_X$, $Y := 3X + U_Y$, $U_X, U_Y \stackrel{iid}{\sim} N(0,1)$

• Update the noise distribution $p_U \rightarrow p_{U\,|\,X=2,Y=6.5}$

$$U_X \sim \delta(2), U_Y \sim \delta(0.5)$$

- Obtain the updated SCM $\mathcal{M}^{X=2,Y=6.5} = (F, p_{U|X=2,Y=6.5})$
- Perform the intervention do(X=1) on $\mathcal{M}^{X=2,Y=6.5}$

$$p^{\mathcal{M}^{X=2,Y=6.5;do(X=1)}}(Y) = \delta(3.5)$$

Recall: $X := U_X$, $Y := 3X + U_Y$, $U_X, U_Y \stackrel{iid}{\sim} N(0,1)$

- Update the noise distribution $p_U \rightarrow p_{U\,|\,X=2,Y=6.5}$

 $U_X \sim \delta(2), U_Y \sim \delta(0.5)$

- Obtain the updated SCM $\mathcal{M}^{X=2,Y=6.5} = (F, p_{U|X=2,Y=6.5})$
- Perform the intervention do(X=1) on $\mathcal{M}^{X=2,Y=6.5}$

$$p^{\mathcal{M}^{X=2,Y=6.5;do(X=1)}}(Y) = \delta(3.5)$$

- Above differs from the interventional distribution $Y \, | \, do(X=1) \sim N(3,1)$



Altitude and Temperature



Altitude and Temperature



- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A)$

Altitude and Temperature



- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A) \label{eq:prod}$
- Entangled factorization

 $p(A,T) = p(T)p(A\,|\,T)$

Altitude and Temperature



Only in the disentangled factorization some components generalize across inteventions/domains

- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A)$
- Entangled factorization

 $p(A,T) = p(T)p(A\,|\,T)$

Altitude and Temperature



Only in the disentangled factorization some components generalize across inteventions/domains

• Austria and Switzerland (CH)

- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A)$
- Entangled factorization

 $p(A,T) = p(T)p(A\,|\,T)$

Altitude and Temperature



- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A)$
- Entangled factorization

 $p(A,T) = p(T)p(A\,|\,T)$

Only in the disentangled factorization some components generalize across inteventions/domains

• Austria and Switzerland (CH)

$$p_{Austria}(A, T) = p_{Austria}(A)p(T \mid A)$$
$$p_{CH}(A, T) = p_{CH}(A)p(T \mid A)$$

Altitude and Temperature



- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A)$
- Entangled factorization $p(A,T) = p(T)p(A\,|\,T) \label{eq:prod}$

Only in the disentangled factorization some components generalize across inteventions/domains

• Austria and Switzerland (CH)

$$p_{Austria}(A, T) = p_{Austria}(A)p(T \mid A)$$
$$p_{CH}(A, T) = p_{CH}(A)p(T \mid A)$$

 $p(T \,|\, A)$ is likely to be the same across countries

Altitude and Temperature



- Disentangled factorization $p(A,T) = p(A)p(T\,|\,A) \label{eq:planck}$
- Entangled factorization $p(A,T) = p(T)p(A\,|\,T)$

Only in the disentangled factorization some components generalize across inteventions/domains

• Austria and Switzerland (CH)

$$p_{Austria}(A,T) = p_{Austria}(A)p(T \mid A)$$
$$p_{CH}(A,T) = p_{CH}(A)p(T \mid A)$$

 $p(T \mid A)$ is likely to be the same across countries p(A) is country-specific

Principle (Independent Causal Mechanisms (ICM))

The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other.

Principle (Independent Causal Mechanisms (ICM))

The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other.

In the two-variable case, say (A,T), it reduces to independence between

- the cause distribution p(A)
- the mechanism producing the effect from the cause $p(T\,|\,A)$

Principle (Independent Causal Mechanisms (ICM))

The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other.

In the two-variable case, say (A,T), it reduces to independence between

- the cause distribution p(A)
- the mechanism producing the effect from the cause $p(T\,|\,A)$

Principle (Sparse Mechanism Shift)

Small distribution changes manifest in a sparse or local way in the causal/disentangled factorization; i.e., they should usually not affect all factors simultaneously.

Principle (Sparse Mechanism Shift)

Small distribution changes manifest in a sparse or local way in the causal/disentangled factorization; i.e., they should usually not affect all factors simultaneously.

In a non-causal/entangled factorization, many terms will be affected simultaneously if we change one of the physical mechanisms responsible for a system's statistical dependencies

Principle (Sparse Mechanism Shift)

Small distribution changes manifest in a sparse or local way in the causal/disentangled factorization; i.e., they should usually not affect all factors simultaneously.

In a non-causal/entangled factorization, many terms will be affected simultaneously if we change one of the physical mechanisms responsible for a system's statistical dependencies

In the Altitude-Temperature setting, if we change country

- we only need to change p(A) if we use the causal factorization
- we need to change both p(T) and $p(A \mid T)$ in the entangled factorization

Table of Contents

1 From Statistical to Causal Models

2 Causal Modeling

3 Causal Discovery

4 Causal Reasoning

5 Conclusions

• Domain knowledge often unavailable or incomplete

- Domain knowledge often unavailable or incomplete
- Need to learn the causal DAG

- Domain knowledge often unavailable or incomplete
- Need to learn the causal DAG Typically using observational (passive) data which are abundant

- Domain knowledge often unavailable or incomplete
- Need to learn the causal DAG Typically using observational (passive) data which are abundant
- Hopeless?

- Domain knowledge often unavailable or incomplete
- Need to learn the causal DAG Typically using observational (passive) data which are abundant
- Hopeless?
- Surprisingly the problem becomes *easier* when the number of variables becomes *higher*

- Domain knowledge often unavailable or incomplete
- Need to learn the causal DAG Typically using observational (passive) data which are abundant
- Hopeless?
- Surprisingly the problem becomes *easier* when the number of variables becomes *higher* because there are nontrivial *conditional independence* properties among the variables implied by a causal structure
Basic idea

- Test which (conditional) independencies can be inferred from the data
- Try to find a graph which implies them

Basic idea

- Test which (conditional) independencies can be inferred from the data
- Try to find a graph which implies them

Assumption (Faithfulness)

The only (conditional) independencies satisfied by $p(\cdot)$ are those implied by the causal Markov condition

Example: faithfulness-SCM

$$X_1 := U_1$$

$$X_2 = \alpha X_1 + U_2$$

$$X_3 = \beta X_1 + \gamma X_2 + U_3$$

with $U_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$.



Causal DAG ${\cal G}$

| - | | | ~ | |
|------|-----|---|----------|--|
| | 110 | | (onconn | |
| C II | | 0 | CONSON | |
| | | | | |

Causal Reasoning

O'Bayes 2025, June 8, Athens

$$X_3 = (\beta + \alpha \gamma)X_1 + \gamma U_2 + U_3$$

$$X_3 = (\beta + \alpha \gamma)X_1 + \gamma U_2 + U_3$$

Thus, if $\beta + \alpha \gamma = 0$, then:

 $X_1 \perp \!\!\!\perp X_3$

$$X_3 = (\beta + \alpha \gamma)X_1 + \gamma U_2 + U_3$$

Thus, if $\beta + \alpha \gamma = 0$, then:

 $X_1 \perp \!\!\!\perp X_3$

But this is **not** implied by the graph G.

$$X_3 = (\beta + \alpha \gamma)X_1 + \gamma U_2 + U_3$$

Thus, if $\beta + \alpha \gamma = 0$, then:

 $X_1 \perp \!\!\!\perp X_3$

But this is **not** implied by the graph G.

Faithfulness is violated

Definition (Markov equivalence)

Two DAGs are said to be Markov equivalent if they encode the same conditional independence (CI) statements.

The set of all DAGs encoding the same CI's is called a Markov equivalence class

Chains, forks and colliders



(a) and (b) imply $X \perp Z \mid Y$ (and no others)

Chains, forks and colliders



(a) and (b) imply $X \perp Z \mid Y$ (and no others)



(c) implies $X \perp \!\!\!\perp Z$ (but $X \not\perp \!\!\!\perp Z \mid Y$)

(a) and (b): same Markov equiv class (c) singleton equivalence class

Markov equivalence: characterization

Result

Two DAG's are Markov equivalent iff they have the same skeleton and the same v-structures

Markov equivalence: characterization

Result

Two DAG's are Markov equivalent iff they have the same skeleton and the same v-structures

Skeleton of Chains (a), Fork (b) and Collider (c)

Markov equivalence: characterization

Result

Two DAG's are Markov equivalent iff they have the same skeleton and the same v-structures

Skeleton of Chains (a), Fork (b) and Collider (c) x - zv-structures

✓ Skeleton estimation

- Test $X \perp \!\!\!\perp Y \mid W$ for all $W \subseteq \mathbf{X} \setminus \{X, Y\}$.
- if no such \boldsymbol{W} is found, connect \boldsymbol{X} and \boldsymbol{Y}
- Expensive

Skeleton estimation

- Test $X \perp \!\!\!\perp Y \mid W$ for all $W \subseteq \mathbf{X} \setminus \{X, Y\}$.
- if no such \boldsymbol{W} is found, connect \boldsymbol{X} and \boldsymbol{Y}
- Expensive
- Edge orientation
 - Direct edges avoiding *v*-structures and cycles.
- PC (Spirtes et al. 2000) more efficient
- **FCI** (handles hidden confounders)

Skeleton estimation

- Test $X \perp \!\!\!\perp Y \mid W$ for all $W \subseteq \mathbf{X} \setminus \{X, Y\}$.
- if no such \boldsymbol{W} is found, connect \boldsymbol{X} and \boldsymbol{Y}
- Expensive
- Edge orientation
 - Direct edges avoiding *v*-structures and cycles.
- PC (Spirtes et al. 2000) more efficient
- **FCI** (handles hidden confounders)
- × Limitations
 - Only a Markov equivalence class is recovered.
 - CI testing is a hard problem

Skeleton estimation

- Test $X \perp \!\!\!\perp Y \mid W$ for all $W \subseteq \mathbf{X} \setminus \{X, Y\}$.
- if no such W is found, connect \boldsymbol{X} and \boldsymbol{Y}
- Expensive
- Edge orientation
 - Direct edges avoiding v-structures and cycles.
- PC (Spirtes et al. 2000) more efficient
- FCI (handles hidden confounders)
- × Limitations
 - Only a Markov equivalence class is recovered.
 - CI testing is a hard problem



- \mathcal{G} : set of DAG's over variables $\{X_1, \ldots, X_n\}$
- $D = {\mathbf{x}_1 \dots, \mathbf{x}_m}$: observed data
- + $S(G \,|\, D)$: score reflecting how well a G -graphical statistical model fits D

- \mathcal{G} : set of DAG's over variables $\{X_1, \ldots, X_n\}$
- $D = {\mathbf{x}_1 \dots, \mathbf{x}_m}$: observed data
- + $S(G \,|\, D)$: score reflecting how well a G -graphical statistical model fits D
- Most methods assume a parametric model which factorises according to ${\cal G}$

- \mathcal{G} : set of DAG's over variables $\{X_1, \ldots, X_n\}$
- $D = {\mathbf{x}_1 \dots, \mathbf{x}_m}$: observed data
- + $S(G \mid D)$: score reflecting how well a G-graphical statistical model fits D
- Most methods assume a parametric model which factorises according to ${\cal G}$
- $S_{BIC}(G \mid D) = \log p(D \mid G, \hat{\theta}^{MLE}) \frac{k}{2} \log m$ k = # of parameters

- \mathcal{G} : set of DAG's over variables $\{X_1, \ldots, X_n\}$
- $D = {\mathbf{x}_1 \dots, \mathbf{x}_m}$: observed data
- + $S(G \mid D)$: score reflecting how well a G-graphical statistical model fits D
- Most methods assume a parametric model which factorises according to ${\cal G}$
- $S_{BIC}(G \mid D) = \log p(D \mid G, \hat{\theta}^{MLE}) \frac{k}{2} \log m$ k = # of parameters
- $S_{BAYES}(G \mid D) = p(D \mid G) = \int_{\Theta} p(D \mid G, \theta) p(\theta \mid G) d\theta$

- \mathcal{G} : set of DAG's over variables $\{X_1, \ldots, X_n\}$
- $D = {\mathbf{x}_1 \dots, \mathbf{x}_m}$: observed data
- + $S(G \,|\, D)$: score reflecting how well a G -graphical statistical model fits D
- Most methods assume a parametric model which factorises according to ${\cal G}$
- $S_{BIC}(G \mid D) = \log p(D \mid G, \hat{\theta}^{MLE}) \frac{k}{2} \log m$ k = # of parameters
- $S_{BAYES}(G \mid D) = p(D \mid G) = \int_{\Theta} p(D \mid G, \theta) p(\theta \mid G) d\theta$ $\hat{G} = \operatorname*{argmax}_{G \in G} S(G \mid D)$

- \mathcal{G} : set of DAG's over variables $\{X_1, \ldots, X_n\}$
- $D = {\mathbf{x}_1 \dots, \mathbf{x}_m}$: observed data
- + $S(G \,|\, D)$: score reflecting how well a G -graphical statistical model fits D
- Most methods assume a parametric model which factorises according to ${\cal G}$
- $S_{BIC}(G \mid D) = \log p(D \mid G, \hat{\theta}^{MLE}) \frac{k}{2} \log m$ k = # of parameters

•
$$S_{BAYES}(G \mid D) = p(D \mid G) = \int_{\Theta} p(D \mid G, \theta) p(\theta \mid G) d\theta$$

 $\hat{G} = \operatorname*{argmax}_{G \in \mathcal{G}} S(G \mid D)$

With a prior p(G) can also use the full posterior

$$p(G \mid D) \propto p(D \mid G)p(G)$$

Table of Contents

1 From Statistical to Causal Models

2 Causal Modeling

3 Causal Discovery

4 Causal Reasoning

5 Conclusions

Two steps

Two steps

(i) *identify* the query, i.e., derive an estimand that only involves *observable* quantities

Two steps

- (i) *identify* the query, i.e., derive an estimand that only involves *observable* quantities
- (ii) make inference on the estimand using data

Outcome \boldsymbol{Y} and binary treatment \boldsymbol{T}

$$\tau := \mathbb{E}[Y \mid do(T=1)] - \mathbb{E}[Y \mid do(T=0)]$$

Outcome \boldsymbol{Y} and binary treatment \boldsymbol{T}

$$\tau := \mathbb{E}[Y \mid do(T=1)] - \mathbb{E}[Y \mid do(T=0)]$$

Outcome Y and continuous X

$$\tau(x') := \left[\frac{\partial}{\partial x} \mathbb{E}[Y \,|\, do(X=x)\right]_{x=x'}$$

Outcome \boldsymbol{Y} and binary treatment \boldsymbol{T}

$$\tau := \mathbb{E}[Y \,|\, do(T=1)] - \mathbb{E}[Y \,|\, do(T=0)]$$

Outcome Y and continuous X

$$\tau(x') := \left[\frac{\partial}{\partial x} \mathbb{E}[Y \,|\, do(X = x)\right]_{x = x'}$$

Treatment effects involve interventional expressions

Outcome \boldsymbol{Y} and binary treatment \boldsymbol{T}

$$\tau := \mathbb{E}[Y \,|\, do(T=1)] - \mathbb{E}[Y \,|\, do(T=0)]$$

Outcome Y and continuous X

$$\tau(x') := \left[\frac{\partial}{\partial x} \mathbb{E}[Y \,|\, do(X = x)\right]_{x=1}$$

Treatment effects involve interventional expressions Causal reasoning answers queries using observational data together with a causal model

Given a causal graph and no hidden confounders

The causal effect can be identified through the interventional distribution

$$p(X_1, \dots, X_n \mid do(X_i = x_i)) = \delta(x_i) \prod_{j \neq i} p(X_j \mid PA_j)$$
(g)

Given a causal graph and no hidden confounders

The causal effect can be identified through the interventional distribution

$$p(X_1, \dots, X_n \mid do(X_i = x_i)) = \delta(x_i) \prod_{j \neq i} p(X_j \mid PA_j)$$
(g)

The interventional distribution of any $X_h \ (h \neq i)$ can be obtained by marginalization

Given a causal graph and no hidden confounders

The causal effect can be identified through the interventional distribution

$$p(X_1, \dots, X_n \mid do(X_i = x_i)) = \delta(x_i) \prod_{j \neq i} p(X_j \mid PA_j)$$
(g)

The interventional distribution of any $X_h \ (h \neq i)$ can be obtained by marginalization

Remarks

Formula (g) has been named

- *g-formula* Robins (1986)
- *truncated factorization* Pearl (2000, 2009)

Given a causal graph and no hidden confounders

The causal effect can be identified through the interventional distribution

$$p(X_1, \dots, X_n \mid do(X_i = x_i)) = \delta(x_i) \prod_{j \neq i} p(X_j \mid PA_j)$$
(g)

The interventional distribution of any X_h $(h \neq i)$ can be obtained by marginalization

Remarks

Formula (g) has been named

- *g-formula* Robins (1986)
- *truncated factorization* Pearl (2000, 2009)

It relies on the independence of causal mechanisms i.e. intervening on a variable leaves the remaining causal conditionals unaffected
Evaluation of treatment effect with three covariates $\{X_1, X_2, X_3\}$



O'Bayes 2025, June 8, Athens

Evaluation of treatment effect with three covariates $\{X_1, X_2, X_3\}$



Factorization of interventional distribution

 $p(y,t,x_1,x_2,x_3 \mid do(T=t)) = \delta(t)p(x_1)p(x_2 \mid x_1)p(y \mid x_2,x_3,t)p(x_3 \mid x_2,t)$

O'Bayes 2025, June 8, Athens

Evaluation of treatment effect with three covariates $\{X_1, X_2, X_3\}$



Factorization of interventional distribution

 $p(y,t,x_1,x_2,x_3 \mid do(T=t)) = \delta(t)p(x_1)p(x_2 \mid x_1)p(y \mid x_2,x_3,t)p(x_3 \mid x_2,t)$

Target distribution p(y | do(T = t))

$$p(y \mid do(T = t)) = \sum_{x_1, x_2, x_3} p(y, t, x_1, x_2, x_3 \mid do(T = t))$$

=
$$\sum_{x_2} \sum_{x_1} p(x_2 \mid x_1) p(x_1) \sum_{x_3} p(y \mid x_2, x_3, t) p(x_3 \mid x_2, t)$$

=
$$\sum_{x_2} p(x_2) p(y \mid x_2, t)$$

O'Bayes 2025, June 8, Athens

$$p(y \mid do(T = t)) = \sum_{x_1, x_2, x_3} p(y, t, x_1, x_2, x_3 \mid do(T = t))$$

=
$$\sum_{x_2} \sum_{x_1} p(x_2 \mid x_1) p(x_1) \sum_{x_3} p(y \mid x_2, x_3, t) p(x_3 \mid x_2, t)$$

=
$$\sum_{x_2} p(x_2) p(y \mid x_2, t)$$

 x_2 is a valid adjustment set

O'Bayes 2025, June 8, Athens

It can be proved using graphical criteria or otherwise that

$$Y \perp X_1 \mid \{T, X_2\}$$
(1.a)
$$X_2 \perp T \mid X_1$$
(1.b)

It can be proved using graphical criteria or otherwise that

$$Y \perp \!\!\perp X_1 \mid \{T, X_2\} \tag{1.a}$$

$$X_2 \bot\!\!\!\bot T \mid X_1 \tag{1.b}$$

$$p(y \mid do(T = t)) = \sum_{x_1, x_2} p(x_1, x_2) p(y \mid x_1, x_2, t), \text{ using (1.a)}$$
(2.a)
$$= \sum_{x_1} p(x_1) \sum_{x_2} p(x_2 \mid x_1, t) p(y \mid x_1, x_2, t), \text{ using (1.b)}$$

$$= \sum_{x_1} p(x_1) p(y \mid x_1, t)$$
(2.b)

It can be proved using graphical criteria or otherwise that

$$Y \perp \!\!\perp X_1 \mid \{T, X_2\} \tag{1.a}$$

$$X_2 \bot\!\!\!\bot T \mid X_1 \tag{1.b}$$

$$p(y \mid do(T = t)) = \sum_{x_1, x_2} p(x_1, x_2) p(y \mid x_1, x_2, t), \text{ using (1.a)}$$
(2.a)
$$= \sum_{x_1} p(x_1) \sum_{x_2} p(x_2 \mid x_1, t) p(y \mid x_1, x_2, t), \text{ using (1.b)}$$

$$= \sum_{x_1} p(x_1) p(y \mid x_1, t)$$
(2.b)

Both $\{x_1, x_2\}$ by (2.a) and $\{x_1\}$ by (2.b) are valid adjustment sets. However $\{x_1, x_3\}$ is not.

Whenever

$$p(y \mid do(T = t)) = \sum_{z} p(z)p(y \mid z, t)$$
(3)

 \boldsymbol{z} is called a valid adjustment set

Whenever

$$p(y \mid do(T = t)) = \sum_{z} p(z)p(y \mid z, t)$$
(3)

 \boldsymbol{z} is called a valid adjustment set

Under causal sufficiency (no hidden variables) there exist graphical criteria to find valid adjustment sets

Whenever

$$p(y \mid do(T = t)) = \sum_{z} p(z)p(y \mid z, t)$$
(3)

 \boldsymbol{z} is called a valid adjustment set

Under causal sufficiency (no hidden variables) there exist graphical criteria to find valid adjustment sets

To estimate the involved quantities in (3) additional assumptions are required in particular *overlap*: for any t and feature values \mathbf{x} , \mathbf{X} , $0 < p(T = t | \mathbf{X} = \mathbf{x}) < 1$

Whenever

$$p(y \mid do(T = t)) = \sum_{z} p(z)p(y \mid z, t)$$
(3)

 \boldsymbol{z} is called a valid adjustment set

Under causal sufficiency (no hidden variables) there exist graphical criteria to find valid adjustment sets

To estimate the involved quantities in (3) additional assumptions are required in particular *overlap*: for any t and feature values \mathbf{x} , \mathbf{X} , $0 < p(T = t | \mathbf{X} = \mathbf{x}) < 1$



Estimation of treatment effect by regression adjustment

Z: adjustment set

Estimation of treatment effect by regression adjustment

Z: adjustment set

Expected value of the outcome \boldsymbol{Y} following an intervention on \boldsymbol{T}

$$\begin{split} \mathbb{E}[Y \mid do(T=t)] &= \sum_{y} yp(y \mid do(T=t)) \\ &= \sum_{y} y \sum_{z} p(z)p(y \mid z, t) \\ &= \sum_{z} p(z) \sum_{y} yp(y \mid z, t) = \sum_{z} p(z)\mathbb{E}[Y \mid z, t] \\ &= \sum_{z} p(z)f(z, t) \end{split}$$

Estimation of treatment effect by regression adjustment

Z: adjustment set

Expected value of the outcome Y following an intervention on T

$$\begin{split} \mathbb{E}[Y \mid do(T=t)] &= \sum_{y} yp(y \mid do(T=t)) \\ &= \sum_{y} y \sum_{z} p(z)p(y \mid z, t) \\ &= \sum_{z} p(z) \sum_{y} yp(y \mid z, t) = \sum_{z} p(z)\mathbb{E}[Y \mid z, t] \\ &= \sum_{z} p(z)f(z, t) \end{split}$$

(Average) Treatment Effect (ATE)

$$\tau = \mathbb{E}[Y \mid do(T=1)] - \mathbb{E}[Y \mid do(T=0)]$$
$$= \sum_{z} p(z)[f(z,1) - f(z,0)]$$

 $\widehat{f}(z,t){:}$ estimator of $\mathbb{E}[Y \mid z,t]$ based on a regression model

 $\widehat{f}(z,t){:}$ estimator of $\mathbb{E}[Y\mid z,t]$ based on a regression model

Regression adjusted plug-in estimator of ATE

$$\hat{\tau}_1 = \frac{1}{m} \sum_{i=1}^m (\hat{f}(z,1) - \hat{f}(z,0))$$

 $\widehat{f}(z,t)$: estimator of $\mathbb{E}[Y \mid z,t]$ based on a regression model

Regression adjusted plug-in estimator of ATE

$$\hat{\tau}_1 = \frac{1}{m} \sum_{i=1}^m (\hat{f}(z,1) - \hat{f}(z,0))$$

An alternative robust estimator

$$\hat{\tau}_2 = \frac{1}{m_1} \sum_{i:t_i=1} (y_i - \hat{f}(z_i, 0)) + \frac{1}{m_0} \sum_{i:t_i=0} (\hat{f}(z_i, 1) - y_i)$$

 m_1 : # obs in the treatment group m_0 # obs in the control group

Further methods to estimate ATE

- Matching and Weighting
- Propensity Score-Methods

Causal inference with unobserved confounders

In general this is not possible

Causal inference with unobserved confounders

In general this is not possible

In some particular situations ATE can still be estimated

In general this is not possible

In some particular situations ATE can still be estimated

• Front-Door Adjustment (Mediator)

In general this is not possible

In some particular situations ATE can still be estimated

- Front-Door Adjustment (Mediator)
- Instrumental Variables (IV)
 - Mendelian randomization: special case

In general this is not possible

In some particular situations ATE can still be estimated

- Front-Door Adjustment (Mediator)
- Instrumental Variables (IV)
 - Mendelian randomization: special case
- Regression Discontinuity Design

DAGs versus Potential Outcomes

Reference Guido Imbens Potential Outcome and Directed Acyclic Graph Approaches to Causality: Relevance for Empirical Practice in Economics. *Journal of Economic Literature*, 2020

DAGs

- J Pearl
 - Precursor: S Wright (path analysis)
 - Computer Science, Statistics, Epidemiology, ML/AI
- Graph captures the way researchers think about causality
- Powerful way to illustrate assumptions
- Systematic way to answer causal queries (do-calculus)
- Useful in complex models (large number of variables)

Potential outcomes

• D Rubin

- Precursors: R Fisher, J Neyman (RCT's)
- Economics, Econometrics, Social sciences
- Critical assumptions (monotonocity, convexity) easier to explain and incorporate
- Connects well to economic theory
- Has established canonical identification strategies for problems with a small number of variables
- Deals nicely with heterogeneity, study designs, estimation

Imbens's Conclusions

• The DAG approach fully deserves the attention of all researchers and users of causal inference.

- The DAG approach fully deserves the attention of all researchers and users of causal inference.
- Two key questions:
 - Should it be the framework of choice for all causal questions, or at least in the social sciences?

- The DAG approach fully deserves the attention of all researchers and users of causal inference.
- Two key questions:
 - Should it be the framework of choice for all causal questions, or at least in the social sciences?
 - Should it be the starting point for teaching about causality?

- The DAG approach fully deserves the attention of all researchers and users of causal inference.
- Two key questions:
 - Should it be the framework of choice for all causal questions, or at least in the social sciences?
 - Should it be the starting point for teaching about causality?
- Imbens's answer to both questions is NO

Table of Contents

1 From Statistical to Causal Models

2 Causal Modeling

3 Causal Discovery

4 Causal Reasoning



• Causality is a key-topic A Gelman & A Vehtari, 2021

What are the Most Important Statistical Ideas of the Past 50 Years?

• Causality is a key-topic A Gelman & A Vehtari, 2021

What are the Most Important Statistical Ideas of the Past 50 Years?

• DAGs or Potential Outcomes Opportunity not a problem Each has its own merits • Causality is a key-topic A Gelman & A Vehtari, 2021

What are the Most Important Statistical Ideas of the Past 50 Years?

- DAGs or Potential Outcomes Opportunity not a problem Each has its own merits
- DAG approach resonates better within the Computer Science community

ML/AI

It features already in several research areas

Promises to have a tremendous impact in the near future
• Causality is a key-topic A Gelman & A Vehtari, 2021

What are the Most Important Statistical Ideas of the Past 50 Years?

- DAGs or Potential Outcomes Opportunity not a problem Each has its own merits
- DAG approach resonates better within the Computer Science community

ML/AI

It features already in several research areas

Promises to have a tremendous impact in the near future

It is a bridge between Statistics-Data Science-AI

• Causal Representation Learning.

 Causal Representation Learning. Learn variables from data

- Causal Representation Learning. Learn variables from data Interventions, reasoning, planning
- Causal Auto encoder

- Causal Representation Learning. Learn variables from data Interventions, reasoning, planning
- Causal Auto encoder
 Build a generative causal model
- Learning transferable mechanisms

- Causal Representation Learning. Learn variables from data Interventions, reasoning, planning
- Causal Auto encoder
 Build a generative causal model
- Learning transferable mechanisms Solving multiple tasks in multiple environments

Selected References

"JUST EXTRAORDINARY." -SCIENCE FRIDAY (NFR)

JUDEA PEARL

AND DANA MACKENZIE

THE BOOK OF

WHY



THE NEW SCIENCE OF CAUSE AND EFFECT



FROM STATISTICAL TO CAUSAL LEARNING

Bernhard Schölkopf Max Planck Institute for Intelligent Systems, Tübingen, Germany bs@tuebingen.mpg.de

Julius von Kügelgen Max Planck Institute for Intelligent Systems, Tübingen, Germany University of Cambridge, United Kingdom jvk8tuebingen. mpg.de

April 4, 2022

ABSTRACT

We describe basic ideas underlying research to build and understand artificially intelligent systems: from symbolic approaches via statistical learning to interventional models relying on concepts of causality. Some of the hard open problems of machine learning and AI are intrinsically related to causality, and progress may require advances in our understanding of how to model and infer causality from data.*

Mathematics Subject Classification 2020

Primary 68T05; Secondary 68Q32, 68T01, 68T10, 68T30, 68T37

Keywords

Causal inference, machine learning, causal representation learning



O'Bayes 2025, June 8, Athens

Causal Reasoning