## Differentiable Causal Discovery with observational and interventional data



Philippe Brouillard\*1





Sébastien Lachapelle\*1



Alexandre Lacoste<sup>2</sup>



Simon Lacoste-Julien<sup>1</sup> \* Equal contribution



Alexandre Drouin<sup>2</sup>

<sup>1</sup> Mila & DIRO, Université de Montréal <sup>2</sup> Element AI / ServiceNow

OATML, University of Oxford February 16, 2021







# Outline

Introduction and Motivation

2 Causal Discovery

Differentiable Causal Discovery with Interventional Data

Conclusions and Future Directions

メロト メタト メヨト メヨト

Introduction and Motivation

・ロト ・四ト ・ヨト ・ヨト

Consider the relationships between altitude (A) and temperature (T)



P(A, T) = P(A|T)P(T)= P(T|A)P(A)

イロト イポト イヨト イヨト

Example taken from [Peters et al., 2017]

## If altitude $\uparrow$ then temperature $\downarrow$

Consider the relationships between altitude (A) and temperature (T)



P(A, T) = P(A|T)P(T)= P(T|A)P(A)

イロト イポト イヨト イヨト

Example taken from [Peters et al., 2017]

#### If altitude $\downarrow$ then temperature $\uparrow$

Consider the relationships between altitude (A) and temperature (T)



Example taken from [Peters et al., 2017]

## Will cooling house #1 make it climb the mountain?

Alexandre Drouin

February 16, 2021 4 / 31

イロト イポト イヨト イヨト

Consider the relationships between altitude (A) and temperature (T)



## Will pushing house #2 down the mountain change its temperature?

Alexandre Drouin

Differentiable Causal Discovery

February 16, 2021 4 / 31

イロト イポト イヨト イヨト

	Overall	Patients with small stones	Patients with large stones
Treatment <i>a</i> : Open surgery	78% (273/350)	<b>93%</b> (81/87)	<b>73%</b> (192/263)
Treatment b: Percutaneous nephrolithotomy	<b>83%</b> (289/350)	87% (234/270)	69% (55/80)

#### Recovery of kidney stone patients

Example taken from Julious & Mullee [1994]

イロト 不良 とくほとくほう

	Overall	Patients with small stones	Patients with large stones
Treatment <i>a</i> : Open surgery	78% (273/350)	<b>93%</b> (81/87)	<b>73%</b> (192/263)
Treatment b: Percutaneous nephrolithotomy	<b>83%</b> (289/350)	87% (234/270)	69% (55/80)

#### Recovery of kidney stone patients

Example taken from Julious & Mullee [1994]

#### • Small stones: treatment a more effective

• Large stones: treatment a more effective

	Overall	Patients with small stones	Patients with large stones
Treatment <i>a</i> : Open surgery	78% (273/350)	<b>93%</b> (81/87)	<b>73%</b> (192/263)
Treatment b: Percutaneous nephrolithotomy	<b>83%</b> (289/350)	87% (234/270)	69% (55/80)

#### Recovery of kidney stone patients

Example taken from Julious & Mullee [1994]

#### • Small stones: treatment a more effective

• Large stones: treatment a more effective

	Overall	Patients with small stones	Patients with large stones
Treatment <i>a</i> : Open surgery	78% (273/350)	<b>93%</b> (81/87)	<b>73%</b> (192/263)
Treatment b: Perçutaneous nephrolithotomy	<b>83%</b> (289/350)	87% (234/270)	69% (55/80)

#### Recovery of kidney stone patients

Example taken from Julious & Mullee [1994]

- Small stones: treatment a more effective
- Large stones: treatment a more effective

All patients: treatment b more effective!

イロト 不得 トイヨト イヨト

# Simpson's paradox: what's really going on?

- $T = \mathsf{Treatment} \in \{A, B\}$
- S =Stone size  $\in$ {small, large}
- $R = Patient recovered \in \{0, 1\}$



< ロ > < 回 > < 回 > < 回 > < 回 >

	Overall	Patients with small stones	Patients with large stones
Treatment a: Open surgery	78% (273/350)	<b>93%</b> (81/87)	<b>73%</b> (192/263)
Treatment b: Percutaneous nephrolithotomy	<b>83%</b> (289/350)	87% (234/270)	69% (55/80)

**Confounding:** patients with small stones are more likely to receive treatment b and also more likely to have a good outcome.

# Simpson's paradox: what's really going on?

- $T = \mathsf{Treatment} \in \{A, B\}$
- S =Stone size  $\in$ {small, large}
- $R = Patient recovered \in \{0, 1\}$



< ロ > < 回 > < 回 > < 回 > < 回 >

	Overall	Patients with small stones	Patients with large stones
Treatment a: Open surgery	78% (273/350)	<b>93%</b> (81/87)	<b>73%</b> (192/263)
Treatment b: Percutaneous nephrolithotomy	<b>83%</b> (289/350)	87% (234/270)	69% (55/80)

**Confounding:** patients with small stones are more likely to receive treatment b and also more likely to have a good outcome.

## Causal vs non-causal questions

• Non-causal: What is the probability of recovery given that the doctor assigned treatment A?

 $\rightarrow P(R = 1 \mid T = A)$ 



• Causal: What's the probability of recovery if treatment A is used irrespective of the size of the stones?

$$\rightarrow P(R = 1 \mid do(T = A))$$



## Causal vs non-causal questions

• Non-causal: What is the probability of recovery given that the doctor assigned treatment A?

 $\mapsto P(R=1 \mid T=A)$ 



• Causal: What's the probability of recovery if treatment A is used irrespective of the size of the stones?

$$\rightarrow P(R = 1 \mid do(T = A))$$



• Random vector 
$$X = (X_1, ..., X_d)$$

- Let G be a directed acyclic graph (DAG)
  - d vertices (one per X<sub>i</sub>)
  - edges indicate causal relationships
- Encodes (conditional) independence constraints (via d-separation, see Koller & Friedman [2009])
- Distribution  $P_X$ :  $p(X) = \prod_{i=1}^d p(X_i | X_{\pi_i^{\mathcal{G}}})$ , where  $\pi_i^{\mathcal{G}}$  = parents of *i* in  $\mathcal{G}$

- Random vector  $X = (X_1, ..., X_d)$
- Let G be a directed acyclic graph (DAG)
  - d vertices (one per X<sub>i</sub>)
  - edges indicate causal relationships
- Encodes (conditional) independence constraints (via d-separation, see Koller & Friedman [2009])
- Distribution  $P_X$ :  $p(X) = \prod_{i=1}^d p(X_i | X_{\pi_i^{\mathcal{G}}})$ , where  $\pi_i^{\mathcal{G}}$  = parents of *i* in  $\mathcal{G}$



- Random vector  $X = (X_1, ..., X_d)$
- Let G be a directed acyclic graph (DAG)
  - d vertices (one per X<sub>i</sub>)
  - edges indicate causal relationships
- Encodes (conditional) independence constraints (via d-separation, see Koller & Friedman [2009])
- Distribution  $P_X$ :  $p(X) = \prod_{i=1}^d p(X_i | X_{\pi_i^{\mathcal{G}}})$ , where  $\pi_i^{\mathcal{G}}$  = parents of *i* in  $\mathcal{G}$



- Random vector  $X = (X_1, ..., X_d)$
- Let G be a directed acyclic graph (DAG)
  - d vertices (one per X<sub>i</sub>)
  - edges indicate causal relationships
- Encodes (conditional) independence constraints (via d-separation, see Koller & Friedman [2009])
- Distribution  $P_X$ :  $p(X) = \prod_{i=1}^d p(X_i | X_{\pi_i^{\mathcal{G}}})$ , where  $\pi_i^{\mathcal{G}}$  = parents of *i* in  $\mathcal{G}$



#### Intervening on the treatment T

- $T = \text{Treatment} \in \{A, B\}$
- S =Stone size  $\in$ {small, large}
- R =Patient recovered  $\in \{0, 1\}$

## Observations



 $p(S)p(T \mid S)p(R \mid S,T)$ 

#### Observational data: may contain biases

$$P(R=1 \mid T=A) \neq P(R=1 \mid do(T=A))$$

イロト イヨト イヨト イヨト

#### Intervening on the treatment T

- $T = \text{Treatment} \in \{A, B\}$
- S =Stone size  $\in$  {small, large}
- R =Patient recovered  $\in \{0, 1\}$

# ObservationsPerfect intervention $\overrightarrow{S}$ $\overrightarrow{S}$ $\overrightarrow{T}$ $\overrightarrow{R}$ $p(S)p(T \mid S)p(R \mid S, T)$ $p(S)\overrightarrow{p}(T)p(R \mid S, T)$

**Perfect intervention:** edges into T are removed (e.g., via randomization) P(R = 1 | T = A) = P(R = 1 | do(T = A))

イロト イロト イヨト イヨト 二日

#### Intervening on the treatment T

- $T = \text{Treatment} \in \{A, B\}$
- S =Stone size  $\in$ {small, large}
- R =Patient recovered  $\in \{0, 1\}$



Imperfect intervention: incoming edges are preserved, conditionals are changed.

イロト イヨト イヨト イヨト

#### Intervening on the treatment T

- $T = \text{Treatment} \in \{A, B\}$
- S =Stone size  $\in \{$ small, large $\}$
- R =Patient recovered  $\in \{0, 1\}$



**A** Important: notice how conditionals that are not under intervention are invariant across distributions (a.k.a, modularity/autonomy).

イロト イヨト イヨト イヨト

• **Objective:** estimate an causal quantity P(R = 1 | do(T = A))

 $\square$  randomization may not be possible (e.g., unethical)

Observational distribution

Interventional distribution

イロト 不得 トイヨト イヨト



Identification: transforming a causal quantity into a purely statistical one

- 🖻 do-calculus [Pearl, 2009] automates this and is even complete [Huang & Valtorta, 2006; Shpitser & Pearl, 2006]
- Example: covariate adjustment

$$P(R = 1 \mid do(T = A)) = \sum_{s \in \{\text{small,large}\}} P(R = 1 \mid T = A, S = s) P(S = s)$$

• **Objective:** estimate an causal quantity P(R = 1 | do(T = A))

 $\mapsto$  randomization may not be possible (e.g., unethical)

Observational distribution

Interventional distribution

イロト 不得 トイヨト イヨト



• Identification: transforming a causal quantity into a purely statistical one

- do-calculus [Pearl, 2009] automates this and is even complete [Huang & Valtorta, 2006; Shpitser & Pearl, 2006]
- Example: covariate adjustment

$$P(R = 1 \mid do(T = A)) = \sum_{s \in \{small, large\}} P(R = 1 \mid T = A, S = s) P(S = s)$$

• **Objective:** estimate an causal quantity P(R = 1 | do(T = A))

 $\square$  randomization may not be possible (e.g., unethical)

Observational distribution

Interventional distribution

イロト 不得 トイヨト イヨト



• Identification: transforming a causal quantity into a purely statistical one

- do-calculus [Pearl, 2009] automates this and is even complete [Huang & Valtorta, 2006; Shpitser & Pearl, 2006]
- Example: covariate adjustment

$$P(R = 1 \mid do(T = A)) = \sum_{s \in \{\text{small,large}\}} P(R = 1 \mid T = A, S = s) P(S = s)$$

• **Objective:** estimate an causal quantity P(R = 1 | do(T = A))

 $\square$  randomization may not be possible (e.g., unethical)

Observational distribution

Interventional distribution

イロト 不得 トイヨト イヨト



• Identification: transforming a causal quantity into a purely statistical one

do-calculus [Pearl, 2009] automates this and is even complete [Huang & Valtorta, 2006; Shpitser & Pearl, 2006]

## What if you don't know the causal graph?

# **Causal Discovery**

▲□▶ ▲圖▶ ▲国▶ ▲国▶

## Problem statement

Observational data

		▲_1	∧ <sub>2</sub>	∧ <sub>3</sub>	
	sample 1	1.2	2.6	0.2	
	sample 2	2.3	5.4	0.5	
	sample n	0.9	1.9	0.1	
_				_	
Int	ervention #1	$X_1 = X_2$	X <sub>3</sub>		
sa	ump <sup>Intervention</sup>	#2 X <sub>1</sub>	X2	X <sub>3</sub>	
sa	um sample <sup>Inte</sup>	rvention #	43 X <sub>1</sub>	X <sub>2</sub>	$X_3$

sample sample 1

sample 2

sample n

V

V V

1.2 2.6 0.2 5.4

2.3

0.9 1.9 0.1

 $X_1$  $X_2$ Algorithm X<sub>3</sub>

イロン イロン イヨン イヨン

Interventional data

sam sample 0.5

#### To make this possible, we need to make assumptions

- Causal sufficiency: no hidden confounding variables
- Markov property: *d-separation* in the graph implies *conditional independence* X<sub>1</sub> ⊥⊥<sub>G</sub> X<sub>2</sub> | Z ⇒ X<sub>1</sub> ⊥⊥<sub>P<sub>X</sub></sub> X<sub>2</sub> | Z
- Faithfulness: conditional independence implies d-separation in the graph

 $X_1 \perp \!\!\!\perp_{P_X} X_2 \mid Z \Longrightarrow X_1 \perp \!\!\!\perp_{\mathcal{G}} X_2 \mid Z$ 

イロン イロン イヨン イヨン

To make this possible, we need to make assumptions

- Causal sufficiency: no hidden confounding variables
- Markov property: *d-separation* in the graph implies *conditional independence*  $X_1 \perp\!\!\!\perp_{\mathcal{G}} X_2 \mid Z \Longrightarrow X_1 \perp\!\!\!\perp_{P_X} X_2 \mid Z$
- Faithfulness: conditional independence implies d-separation in the graph

 $X_1 \perp \!\!\!\perp_{P_X} X_2 \mid Z \Longrightarrow X_1 \perp \!\!\!\perp_{\mathcal{G}} X_2 \mid Z$ 

イロト イヨト イヨト イヨト

To make this possible, we need to make assumptions

- Causal sufficiency: no hidden confounding variables
- Markov property: *d-separation* in the graph implies *conditional independence*  $X_1 \perp\!\!\!\perp_G X_2 \mid Z \Longrightarrow X_1 \perp\!\!\!\!\perp_{P_X} X_2 \mid Z$
- Faithfulness: conditional independence implies d-separation in the graph

 $X_1 \perp \!\!\!\perp_{P_X} X_2 \mid Z \Longrightarrow X_1 \perp \!\!\!\perp_{\mathcal{G}} X_2 \mid Z$ 

・ロト ・ 日 ト ・ 日 ト ・ 日 ト ・

To make this possible, we need to make assumptions

- Causal sufficiency: no hidden confounding variables
- Markov property: *d-separation* in the graph implies *conditional independence*  $X_1 \perp\!\!\!\perp_G X_2 \mid Z \Longrightarrow X_1 \perp\!\!\!\!\perp_{P_X} X_2 \mid Z$
- Faithfulness: conditional independence implies d-separation in the graph

 $X_1 \perp\!\!\!\perp_{P_X} X_2 \,|\, Z \Longrightarrow X_1 \perp\!\!\!\!\perp_{\mathcal{G}} X_2 \,|\, Z$ 

ヘロト 人間 ト 人 ヨト 人 ヨト

To make this possible, we need to make assumptions

- Causal sufficiency: no hidden confounding variables
- Markov property: *d-separation* in the graph implies *conditional independence*  $X_1 \perp\!\!\!\perp_G X_2 \mid Z \Longrightarrow X_1 \perp\!\!\!\!\perp_{P_X} X_2 \mid Z$

• Faithfulness: conditional independence implies d-separation in the graph

 $X_1 \perp\!\!\!\perp_{P_X} X_2 \,|\, Z \Longrightarrow X_1 \perp\!\!\!\!\perp_{\mathcal{G}} X_2 \,|\, Z$ 

These last two assumptions guarantee an equivalence between properties of the data and properties of the graph

イロト 不得 トイヨト イヨト

## Score-based causal discovery

- Idea: find the DAG that maximizes a score function (S)
  - E.g., data likelihood + sparsity prior
  - Consistency: need to demonstrate that the score leads to the true solution



Problem: search space grows superexponentially with variables



• Examples: Greedy Equivalence Search [Chickering, 2003], DAG with NO TEARS [Zheng et al., 2018]

イロト イポト イヨト イヨト

## Score-based causal discovery

- Idea: find the DAG that maximizes a score function (S)
  - E.g., data likelihood + sparsity prior
  - Consistency: need to demonstrate that the score leads to the true solution

$$\hat{\mathcal{G}} \in rg\max_{\mathcal{G} \in \mathit{DAG}} \mathcal{S}(\mathcal{G})$$

• Problem: search space grows superexponentially with variables



Examples: Greedy Equivalence Search [Chickering, 2003], DAG with NO TEARS [Zheng et al., 2018]

イロト イポト イヨト イヨト
# Score-based causal discovery

- Idea: find the DAG that maximizes a score function (S)
  - E.g., data likelihood + sparsity prior
  - Consistency: need to demonstrate that the score leads to the true solution

$$\hat{\mathcal{G}} \in rg\max_{\mathcal{G} \in \mathit{DAG}} \mathcal{S}(\mathcal{G})$$

• Problem: search space grows superexponentially with variables



• Examples: Greedy Equivalence Search [Chickering, 2003], DAG with NO TEARS [Zheng et al., 2018]

イロト イボト イヨト イヨト

# Identifiability of the DAG



Without making more assumptions, observational data only allows identification up to a Markov equivalence class (MEC) [Verma & Pearl, 1991]

< ロ > < 回 > < 回 > < 回 > < 回 >

## Can you shrink the equivalence class?



Interventional Markov equivalence classes

• An I-MEC is a subset of the MEC (see Eberhardt et al. [2005])

• Example: gene knockout/knockdown experiments in biology [Divit et al., 2016]

イロン イロン イヨン イヨン

## Can you shrink the equivalence class?



Interventional Markov equivalence classes

- An I-MEC is a subset of the MEC (see Eberhardt et al. [2005])
- Example: gene knockout/knockdown experiments in biology [Dixit et al., 2016]

イロン イロン イヨン イヨン

# Can you shrink the equivalence class?



Interventional Markov equivalence classes

- An I-MEC is a subset of the MEC (see Eberhardt et al. [2005])
- Example: gene knockout/knockdown experiments in biology [Dixit et al., 2016]

< ロ > < 回 > < 回 > < 回 > < 回 >

Differentiable Causal Discovery with Interventional Data

イロト 不良 とくほとくほう

### Differentiable Causal Discovery from Interventional Data

Philippe Brouillard\* Mila, Université de Montréal Sébastien Lachapelle\* Mila, Université de Montréal Alexandre Lacoste Element AI

Simon Lacoste-Julien Mila, Université de Montréal Canada CIFAR AI Chair Alexandre Drouin Element AI

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

#### **DCDI Fact Sheet:**

- Type: score-based
- Search strategy: continuous-constrained optimization [Zheng et al., 2018]
- Data: observational and interventional (perfect/imperfect)
- Theoretical guarantees: consistency in the limit of infinite data

イロト 不得 トイヨト イヨト



 $p^{(1)}(x_1,...,x_4) := p^{(1)}(x_1)p^{(1)}(x_3|x_1)p^{(1)}(x_2|x_1,x_3)p^{(1)}(x_4|x_1,x_2,x_3)$ 

イロト イヨト イヨト イヨト 二日



 $p^{(2)}(x_1,...,x_4) := p^{(1)}(x_1)p^{(1)}(x_3|x_1)p^{(2)}(x_2|x_1,x_3)p^{(1)}(x_4|x_1,x_2,x_3)$ 

Alexandre Drouin

February 16, 2021 19 / 31

イロト イヨト イヨト イヨト 二日



 $p^{(3)}(x_1,...,x_4) := p^{(3)}(x_1)p^{(1)}(x_3|x_1)p^{(1)}(x_2|x_1,x_3)p^{(3)}(x_4|x_1,x_2,x_3)$ 

Alexandre Drouin

February 16, 2021 19 / 31

イロト イヨト イヨト イヨト 二日



イロト イヨト イヨト イヨト

# DCDI: problem setting and notation

• We observe *d* variables which are *causally sufficient*, i.e. no hidden confounders.

Causal DAG = 
$$\underbrace{G}_{Adjacency matrix} \begin{pmatrix} 0 & 0 & 1\\ 1 & 0 & 1\\ 0 & 0 & 0 \end{pmatrix} \in \{0, 1\}^{d \times d}$$

• We have *K*, potentially **imperfect**, interventions which can target multiple variables simultaneously.

$$\underbrace{I} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix} \in \{0, 1\}^{K \times d}$$

Intervention matrix

・ロト ・四ト ・ヨト ・ヨト

# DCDI: problem setting and notation

• We observe *d* variables which are *causally sufficient*, i.e. no hidden confounders.

Causal DAG = 
$$\underbrace{\begin{array}{c} \mathbf{G} = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}}_{\text{Adjacency matrix}} \in \{0, 1\}^{d \times d}$$

• We have *K*, potentially **imperfect**, interventions which can target multiple variables simultaneously.

$$\underbrace{I = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix}}_{K \times a} \in \{0, 1\}^{K \times a}$$

Intervention matrix

・ロト ・回ト ・ヨト ・ヨト

# DCDI: problem setting and notation

• We observe *d* variables which are *causally sufficient*, i.e. no hidden confounders.

Causal DAG = 
$$\underbrace{\begin{array}{c} G \\ G \end{array}}_{Adjacency matrix} \left( \begin{array}{c} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 0 \end{array} \right) \in \{0, 1\}^{d \times d}$$

• We have *K*, potentially **imperfect**, interventions which can target multiple variables simultaneously.

$$\underbrace{I}_{\text{Intervention matrix}}^{(0)} \in \{0,1\}^{K \times d}$$

イロン イ団 とく ヨン イヨン



The graph adjacency matrix acts as a mask that filters the input variables.

	adro.	1 Irou	
I GAGH	IUI C		

イロト イヨト イヨト イヨト



Each conditional distribution is estimated by a distinct neural network.

Alexandre Drouin			
Alexandre Drouin			 
	нхани	11111	

イロト イヨト イヨト イヨト



The joint likelihood is calculated as the product of conditional distributions (obs/int)

		J F	<b>`</b>	
- M	exant	ле с		

February 16, 2021 21 / 31

イロト 不得 トイヨト イヨト



The intervention matrix *I* activates the right set of parameters.

		-	
	dro	Drou	
			_

Differentiable Causal Discovery

February 16, 2021 21 / 31

イロト イポト イヨト イヨト

# DCDI: graph scoring function (discrete)

• We suggest maximizing this score over the space of DAGs:

$$\mathcal{S}_{I^{\bullet}}(\mathbf{G}) := \sup_{\phi} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \mathbf{G}, \mathbf{I^{*}}, \phi) - \lambda \|\mathbf{G}\|_{0}$$
  
Sparsity regularization Sparsity regularization

• Search: discrete search over DAGs  $\rightarrow$  continuous-constrained opt. problem [Zheng et al., 2018]

< ロ > < 回 > < 回 > < 回 > < 回 >

# DCDI: graph scoring function (discrete)

• We suggest maximizing this score over the space of DAGs:

$$\mathcal{S}_{I^{\ast}}(\mathbf{G}) := \sup_{\phi} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \mathbf{G}, \mathbf{I^{\ast}}, \phi) - \lambda \|\mathbf{G}\|_{0}$$
  
Sparsity regularization Sparsity regularization

• Search: discrete search over DAGs  $\rightarrow$  continuous-constrained opt. problem [Zheng et al., 2018]

イロト イヨト イヨト イヨト

# DCDI: theoretical justification

- $\mathcal{G}^* = \text{ground-truth DAG}$
- $I^* =$  ground-truth intervention matrix

 $\hat{\mathcal{G}} \in \operatorname{arg\,max}_{\mathcal{G} \in \mathsf{DAG}} \mathcal{S}_{I^*}(\mathcal{G})$  is the estimator.

#### I heorem (Identification via score maximization)

Suppose  $I_{1,:}^* = \emptyset$ . Given that

- Each variable is individually targeted by an intervention;
- In the model has enough capacity to express the ground truth;
- ) The regularization coefficient  $\lambda > 0$  is small enough;
- And some more technical assumptions, e.g. 1\*-faithfulness... (See paper, then

 $\hat{\mathcal{G}} = \mathcal{G}^*$ 

#### More general result

Without the first assumption, we can identify the  $I^*$ -Markov equivalence class<sup>a</sup> of  $\mathcal{G}^*$ .

'We use the notion of I\*-Markov equivalence of Yang et al. [2018].

# DCDI: theoretical justification

- $\mathcal{G}^* = \text{ground-truth DAG}$
- $I^* =$  ground-truth intervention matrix

 $\hat{\mathcal{G}} \in \operatorname{arg\,max}_{\mathcal{G} \in \mathsf{DAG}} \mathcal{S}_{I^*}(\mathcal{G})$  is the estimator.

### Theorem (Identification via score maximization)

Suppose  $I_{1,:}^* = \emptyset$ . Given that

- Each variable is individually targeted by an intervention;
- One of the second truth of the second truth
- The regularization coefficient  $\lambda > 0$  is small enough;
- And some more technical assumptions, e.g. I\*-faithfulness... (See paper)

then

$$\hat{\mathcal{G}} = \mathcal{G}^*$$

#### More general result

Without the first assumption, we can identify the  $I^*$ -Markov equivalence class<sup>a</sup> of  $\mathcal{G}^*$ .

'We use the notion of I\*-Markov equivalence of Yang et al. [2018].

# DCDI: theoretical justification

- $\bullet \ \mathcal{G}^* = \mathsf{ground-truth} \ \mathsf{DAG}$
- $I^* =$  ground-truth intervention matrix

 $\hat{\mathcal{G}} \in \operatorname{arg} \max_{\mathcal{G} \in \mathsf{DAG}} \mathcal{S}_{I^*}(\mathcal{G})$  is the estimator.

## Theorem (Identification via score maximization)

Suppose  $I_{1,:}^* = \emptyset$ . Given that

- Each variable is individually targeted by an intervention;
- One of the second truth of the second truth
- The regularization coefficient  $\lambda > 0$  is small enough;
- And some more technical assumptions, e.g. I\*-faithfulness... (See paper)

then

$$\hat{\mathcal{G}}=\mathcal{G}^*$$

### More general result

Without the first assumption, we can identify the  $I^*$ -Markov equivalence class<sup>a</sup> of  $\mathcal{G}^*$ .

<sup>a</sup>We use the notion of *I*\*-Markov equivalence of Yang et al. [2018].

Alexandre Drouin

$$\mathcal{S}_{I^{\bullet}}(\boldsymbol{G}) := \sup_{\phi} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \boldsymbol{G}, \boldsymbol{I^{*}}, \phi) - \lambda \|\boldsymbol{G}\|_{0}$$

Relaxation where  $G_{ij} \sim \text{Bernoulli}(\sigma(\Lambda_{ij}))$ , with  $\sigma(\cdot) :=$  sigmoid function

$$\hat{\mathcal{S}}_{I}(\Lambda) := \sup_{\phi} \mathbb{E}_{\mathbf{G} \sim \sigma(\Lambda)} \left[ \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \mathbf{G}, \mathbf{I}^{*}, \phi) - \lambda ||\mathbf{G}||_{0} \right]$$

Optimize for  $\Lambda$  under acyclicity constraint

< ロ > < 回 > < 回 > < 回 > < 回 >

$$\sup_{\Delta} \hat{S}_{I^{\bullet}}(\Lambda) \quad s.t. \quad \underbrace{\operatorname{Tr}\left(e^{\sigma(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$$

$$[Zheng et al., 2018]$$

$$\hat{\mathcal{S}}_{I^{\bullet}}(\Lambda) := \sup_{\phi} \mathop{\mathbb{E}}_{\boldsymbol{G} \sim \sigma(\Lambda)} \left[ \sum_{k=1}^{K} \mathop{\mathbb{E}}_{X \sim p^{(k)}} \log f^{(k)}(X; \boldsymbol{G}, \boldsymbol{I}^{\bullet}, \phi) - \lambda ||\boldsymbol{G}||_{0} \right]$$

Optimize for  $\Lambda$  under acyclicity constraint

メロト メタト メヨト メヨト

$$\sup_{\Lambda} \hat{S}_{I^{*}}(\Lambda) \quad s.t. \quad \underbrace{\operatorname{Tr}\left(e^{\sigma(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$$
[Zheng et al., 2018]

17

Optimize for  $\Lambda$  under acyclicity constraint

メロト メタト メヨト メヨト

$$\sup_{\Lambda} \hat{S}_{I^{\bullet}}(\Lambda) \quad s.t. \quad \underbrace{\operatorname{Tr}\left(e^{\sigma(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$$
[Zheng et al., 2018]

 $\hat{\mathcal{S}}_{I^*}(\Lambda) := \sup_{\phi} \mathop{\mathbb{E}}_{\substack{G \sim \sigma(\Lambda)}} \left[ \sum_{k=1}^{K} \mathop{\mathbb{E}}_{X \sim p^{(k)}} \log f^{(k)}(X; \mathbf{G}, \mathbf{I}^*, \phi) - \lambda ||\mathbf{G}||_0 \right]$ 

Optimize for  $\Lambda$  under acyclicity constraint

イロン イヨン イヨン イヨン 三日

$$\sup_{\Delta} \hat{\mathcal{S}}_{I^{*}}(\Delta) \quad s.t. \quad \underbrace{\operatorname{Tr}\left(e^{\sigma(\Delta)}\right) - d = 0}_{\text{Acyclicity constraint}}$$

$$[Zheng et al., 2018]$$

$$\mathcal{S}_{I^*}(\mathbf{G}) := \sup_{\phi} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \mathbf{G}, \mathbf{I}^*, \phi) - \lambda \|\mathbf{G}\|_0$$

Relaxation where  $G_{ij} \sim \text{Bernoulli}(\sigma(\Lambda_{ij}))$ , with  $\sigma(\cdot) :=$  sigmoid function

$$\hat{\mathcal{S}}_{I^*}(\Lambda) := \sup_{\phi} \mathbb{E}_{\substack{\boldsymbol{G} \sim \sigma(\Lambda)}} \left[ \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \boldsymbol{G}, \boldsymbol{I^*}, \phi) - \lambda ||\boldsymbol{G}||_0 \right]$$

Optimize for  $\Lambda$  under acyclicity constraint

イロト イヨト イヨト イヨト

 $\sup_{\Lambda} \hat{\mathcal{S}}_{I^*}(\Lambda) \quad s.t. \quad \underbrace{\operatorname{Tr}\left(e^{\sigma(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$ [Zheng et al., 2018]

• Optimize jointly 
$$\Lambda$$
 and  $\phi$  (NN parameters)  

$$\max_{\phi,\Lambda} \mathop{\mathbb{E}}_{G\sim\sigma(\Lambda)} \left[ \sum_{k=1}^{K} \mathop{\mathbb{E}}_{X\sim p^{(k)}} \log f^{(k)}(X; G, I^*, \phi) - \lambda ||G||_0 \right] \text{ s.t. } \underbrace{\operatorname{Tr}\left(e^{\sigma(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$$

- Optimization: Augmented Lagrangian Method + RMSprop (as in Lachapelle et al. [2020])
- Discrete sampling: gradient w.r.t. A estimated via Gumbel-Softmax Straight-Through estimator [Jang et al., 2017].
- Masks and/or the Gumbel-Softmax estimator were used before in causal discovery e.g., Kalainathan et al. [2018]; Ng et al. [2019]; Bengio et al. [2019]; Ke et al. [2019]

イロト イポト イヨト イヨト

• Optimize jointly  $\Lambda$  and  $\phi$  (NN parameters)

$$\max_{\phi, \Lambda} \mathbb{E}_{\mathbf{G} \sim \boldsymbol{\sigma}(\Lambda)} \left[ \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X; \mathbf{G}, \mathbf{I}^{*}, \phi) - \lambda ||\mathbf{G}||_{0} \right] \text{ s.t. } \underbrace{\operatorname{Tr}\left(e^{\boldsymbol{\sigma}(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$$

- Optimization: Augmented Lagrangian Method + RMSprop (as in Lachapelle et al. [2020])
- Discrete sampling: gradient w.r.t. A estimated via Gumbel-Softmax Straight-Through estimator [Jang et al., 2017].
- Masks and/or the Gumbel-Softmax estimator were used before in causal discovery e.g., Kalainathan et al. [2018]; Ng et al. [2019]; Bengio et al. [2019]; Ke et al. [2019]

イロト イボト イヨト イヨト

• Optimize jointly  $\Lambda$  and  $\phi$  (NN parameters)

$$\max_{\phi, \Lambda} \mathop{\mathbb{E}}_{\boldsymbol{G} \sim \boldsymbol{\sigma}(\Lambda)} \left[ \sum_{k=1}^{K} \mathop{\mathbb{E}}_{\boldsymbol{X} \sim p^{(k)}} \log f^{(k)}(\boldsymbol{X}; \boldsymbol{G}, \boldsymbol{I}^{*}, \phi) - \lambda ||\boldsymbol{G}||_{0} \right] \text{ s.t. } \underbrace{\operatorname{Tr}\left(e^{\boldsymbol{\sigma}(\Lambda)}\right) - d = 0}_{\text{Acyclicity constraint}}$$

- Optimization: Augmented Lagrangian Method + RMSprop (as in Lachapelle et al. [2020])
- Discrete sampling: gradient w.r.t. A estimated via Gumbel-Softmax Straight-Through estimator [Jang et al., 2017; Maddison et al., 2017].
- Masks and/or the Gumbel-Softmax estimator were used before in causal discovery e.g., Kalainathan et al. [2018]; Ng et al. [2019]; Bengio et al. [2019]; Ke et al. [2019]

イロト 不得 トイヨト イヨト

• Optimize jointly  $\Lambda$  and  $\phi$  (NN parameters)

$$\max_{\phi, \Lambda} \mathop{\mathbb{E}}_{\boldsymbol{G} \sim \boldsymbol{\sigma}(\Lambda)} \left[ \sum_{k=1}^{K} \mathop{\mathbb{E}}_{X \sim p^{(k)}} \log f^{(k)}(X; \boldsymbol{G}, \boldsymbol{I}^{*}, \phi) - \lambda ||\boldsymbol{G}||_{0} \right] \text{ s.t. } \underbrace{\operatorname{Tr}\left(e^{\boldsymbol{\sigma}(\Lambda)}\right) - d = 0}_{\operatorname{Acyclicity constraint}}$$

- Optimization: Augmented Lagrangian Method + RMSprop (as in Lachapelle et al. [2020])
- Discrete sampling: gradient w.r.t. A estimated via Gumbel-Softmax Straight-Through estimator [Jang et al., 2017; Maddison et al., 2017].
- Masks and/or the Gumbel-Softmax estimator were used before in causal discovery e.g., Kalainathan et al. [2018]; Ng et al. [2019]; Bengio et al. [2019]; Ke et al. [2019]

イロト 不得 トイヨト イヨト

# Result: structure learning via continuous optimization



Optimizing the objective gradually prunes anti-causal edges from the graph

イロト イヨト イヨト イヨト

## DCDI: interventions with unknown targets

• Until now we assumed that 1\* was known, i.e. we knew which variables were targeted.

• What if we don't? (e.g., as in Ke et al. [2019])

$$\begin{split} \mathcal{S}(\pmb{G},\pmb{I}) &:= \sup_{\boldsymbol{\phi}} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X;\pmb{G},\pmb{I},\boldsymbol{\phi}) - \lambda \|\pmb{G}\|_{0} - \lambda_{I} \|\pmb{I}\|_{0} \\ &\text{Intervention matrix} \\ &\text{is learned} \\ \end{split}$$

• Learning:

• Can do the same relaxation  $I_{kj} \sim \text{Bernoulli}(\sigma(\beta_{kj}))$ .

• Optimize jointly for  $\phi$ ,  $\Lambda$  and  $\beta$ .

• Theory: We showed the same guarantee holds for this score!

		-		

イロン イロン イヨン イヨン

## DCDI: interventions with unknown targets

- Until now we assumed that 1\* was known, i.e. we knew which variables were targeted.
- What if we don't? (e.g., as in Ke et al. [2019])

$$\begin{split} \mathcal{S}(\pmb{G},\pmb{I}) &:= \sup_{\substack{\phi \\ \text{is learned}}} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X;\pmb{G},\pmb{I},\phi) - \lambda \|\pmb{G}\|_{0} - \lambda_{I} \|\pmb{I}\|_{0} \\ & \text{Additional sparsity regularizer} \end{split}$$

• Learning:

• Can do the same relaxation  $I_{kj} \sim \text{Bernoulli}(\sigma(\beta_{kj}))$ .

• Optimize jointly for  $\phi$ ,  $\Lambda$  and  $\beta$ .

• Theory: We showed the same guarantee holds for this score!

イロン イロン イヨン イヨン

## DCDI: interventions with unknown targets

- Until now we assumed that  $I^*$  was known, i.e. we knew which variables were targeted.
- What if we don't? (e.g., as in Ke et al. [2019]) Learn it!

$$\begin{split} \mathcal{S}(\textbf{G},\textbf{I}) &:= \sup_{\substack{\phi \\ \text{is learned}}} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X;\textbf{G},\textbf{I},\phi) - \lambda \|\textbf{G}\|_{0} - \underbrace{\lambda_{I} \|\textbf{I}\|_{0}}_{\text{regularizer}} \\ & \text{Additional sparsity} \\ \text{regularizer} \end{split}$$

• Learning:

Can do the same relaxation  $I_{kj} \sim \text{Bernoulli}(\sigma(\beta_{kj}))$ 

• Optimize jointly for  $\phi$ ,  $\Lambda$  and  $\beta$ .

Theory: We showed the same guarantee holds for this score!

イロト イヨト イヨト イヨト
#### DCDI: interventions with unknown targets

- Until now we assumed that  $I^*$  was known, i.e. we knew which variables were targeted.
- What if we don't? (e.g., as in Ke et al. [2019]) Learn it!

$$\begin{split} \mathcal{S}(\pmb{G},\pmb{I}) &:= \sup_{\substack{\phi \\ \text{is learned}}} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X;\pmb{G},\pmb{I},\phi) - \lambda \|\pmb{G}\|_{0} - \underbrace{\lambda_{I} \|\pmb{I}\|_{0}}_{\text{Additional sparsity regularizer}} \end{split}$$

• Learning:

• Can do the same relaxation  $I_{kj} \sim \text{Bernoulli}(\sigma(\beta_{kj}))$ .

• Optimize jointly for  $\phi$ ,  $\Lambda$  and  $\beta$ .

Theory: We showed the same guarantee holds for this score!

#### DCDI: interventions with unknown targets

- Until now we assumed that 1\* was known, i.e. we knew which variables were targeted.
- What if we don't? (e.g., as in Ke et al. [2019]) Learn it!

$$\begin{split} \mathcal{S}(\pmb{G},\pmb{I}) &:= \sup_{\substack{\phi \\ \text{is learned}}} \sum_{k=1}^{K} \mathbb{E}_{X \sim p^{(k)}} \log f^{(k)}(X;\pmb{G},\pmb{I},\phi) - \lambda \|\pmb{G}\|_{0} - \underbrace{\lambda_{I} \|\pmb{I}\|_{0}}_{\text{Additional sparsity regularizer}} \end{split}$$

• Learning:

• Can do the same relaxation  $I_{kj} \sim \text{Bernoulli}(\sigma(\beta_{kj}))$ .

- Optimize jointly for  $\phi$ ,  $\Lambda$  and  $\beta$ .
- Theory: We showed the same guarantee holds for this score!

イロト 不得 トイヨト イヨト

## DCDI: choice of density function $\tilde{f}$



# • Gaussian: corresponds to a non-linear + additive noise assumption on the functional form of causal mechanisms

- Identification guaranteed if it actually holds in the distribution [Peters et al., 2014]
- Deep sigmoidal flow: a type of normalizing flow that was shown to be a universal density approximator [Huang et al., 2018]
  - No assumption on functional forms
  - Identification guaranteed by our Thm 1 (with enough interventions)

## DCDI: choice of density function $\tilde{f}$



- Gaussian: corresponds to a non-linear + additive noise assumption on the functional form of causal mechanisms
  - Identification guaranteed if it actually holds in the distribution [Peters et al., 2014]
- Deep sigmoidal flow: a type of normalizing flow that was shown to be a universal density approximator [Huang et al., 2018]
  - No assumption on functional forms
  - Identification guaranteed by our Thm 1 (with enough interventions)

### Results - Structural Hamming Distance (lower is better)

 $\label{eq:constant} \begin{array}{l} \textbf{DCDI-G} = \text{DCDI} \text{ with Gaussian density} \\ \textbf{DCDI-DSF} = \text{DCDI} \text{ with deep sigmoidal flow} \end{array}$ 

ANM = nonlinear with additive noise NN = nonlinear (no additive noise) e = average number of parents



Figure: Known target interventions (20 nodes)



Figure: Unknown target interventions (20 nodes)

< ロ > < 回 > < 回 > < 回 > < 回 >

### Results - Structural Hamming Distance (lower is better)

DCDI-G = DCDI with Gaussian density DCDI-DSF = DCDI with deep sigmoidal flow  $\begin{array}{l} \textbf{ANM} = \text{nonlinear with additive noise} \\ \textbf{NN} = \text{nonlinear (no additive noise)} \\ \textbf{e} = \text{average number of parents} \end{array}$ 



Figure: Known target interventions (20 nodes)



Figure: Unknown target interventions (20 nodes)

< ロ > < 回 > < 回 > < 回 > < 回 >

### Results - Structural Hamming Distance (lower is better)

$$\label{eq:constant} \begin{split} \textbf{DCDI-G} &= \text{DCDI} \text{ with Gaussian density} \\ \textbf{DCDI-DSF} &= \text{DCDI} \text{ with deep sigmoidal flow} \end{split}$$

ANM = nonlinear with additive noise NN = nonlinear (no additive noise) e = average number of parents



Figure: Known target interventions (20 nodes)



Figure: Unknown target interventions (20 nodes)

## Conclusions and Future Directions

#### We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

< ロ > < 回 > < 回 > < 回 > < 回 >

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

#### Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

#### Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

We proposed DCDI, a causal discovery algorithm that:

- is theoretically grounded
- supports perfect, imperfect and unknown-target interventions
- scales well with sample size (compared to methods using kernel-based independence tests)
- achieves state-of-the-art performance, especially on denser graphs

Future work:

- More extensive evaluation: beyond synthetic data, violate assumptions [Gentzel et al., 2019]
- Relax causal sufficiency: allow for hidden confounders [Bhattacharya et al., 2020]
- Time series: non-stationnarities as imperfect interventions [Pamfil et al., 2020]
- Learning variable representations: not being agnostic to the nature of variables

## Thank you!





Philippe Brouillard \*1

Sébastien Lachapelle\*1



Alexandre Lacoste<sup>2</sup>



Simon Lacoste-Julien<sup>1</sup>



Alexandre Drouin<sup>2</sup>

<sup>1</sup> Mila & DIRO, Université de Montréal

<sup>2</sup> Element AI / ServiceNow \* Equal contribution

イロト 不良 とくほとくほう

#### O https://github.com/slachapelle/dcdi

alexandre.drouin@servicenow.com

♥ @alexandredrouin

#### References

- Bengio, Y., Deleu, T., Rahaman, N., Ke, R., Lachapelle, S., Bilaniuk, O., Goyal, A., & Pal, C. (2019). A meta-transfer objective for learning to disentangle causal mechanisms. arXiv preprint arXiv:1901.10912.
- Bhattacharya, R., Nagarajan, T., Malinsky, D., & Shpitser, I. (2020). Differentiable causal discovery under unmeasured confounding. arXiv preprint arXiv:2010.06978.

Chickering, D. (2003). Optimal structure identification with greedy search. Journal of Machine Learning Research.

- Dixit, A., Parnas, O., Li, B., Chen, J., Fulco, C. P., Jerby-Arnon, L., Marjanovic, N. D., Dionne, D., Burks, T., Raychowdhury, R., Adamson, B., Norman, T. M., Lander, E. S., Weissman, J. S., Friedman, N., & Regev, A. (2016). Perturb-seq: Dissecting molecular circuits with scalable single-cell ma profiling of pooled genetic screens. Cell, 167(7), 1853–1866.e17.
- Eberhardt, F., Glymour, C., & Scheines, R. (2005). On the number of experiments sufficient and in the worst case necessary to identify all causal relations among n variables. In Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence, UAI'05 (pp. 178184). Artington, Virginia, USA: AUAI Press.
- Gentzel, A., Garant, D., & Jensen, D. (2019). The case for evaluating causal models using interventional measures and empirical data. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox, & R. Gamett (Eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019. December 8-14, 2019. Vancouver, BC, Canada (pp. 11717–11727).

Huang, C.-W., Krueger, D., Lacoste, A., & Courville, A. (2018). Neural autoregressive flows.

- Huang, Y. & Valtorta, M. (2006). Pearl's calculus of intervention is complete. In Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence, UAI'06 (pp. 217224). Arlington, Virginia, USA: AUAI Press.
- Jang, E., Gu, S., & Poole, B. (2017). Categorical reparameterization with gumbel-softmax. Proceedings of the 34th International Conference on Machine Learning.
- Julious, S. A. & Mullee, M. A. (1994). Confounding and simpson's paradox. Bmj, 309(6967), 1480-1481.
- Kalainathan, D., Goudet, O., Guyon, I., Lopez-Paz, D., & Sebag, M. (2018). Sam: Structural agnostic model, causal discovery and penalized adversarial learning. arXiv preprint arXiv:1803.04929.
- Ke, N. R., Bilaniuk, O., Goyal, A., Bauer, S., Larochelle, H., Pal, C., & Bengio, Y. (2019). Learning neural causal models from unknown interventions. arXiv preprint arXiv:1910.01075. Koller, D. & Friedman, N. (2009). Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning. MIT Press.
- Lachapelle, S., Brouillard, P., Deleu, T., & Lacoste-Julien, S. (2020). Gradient-based neural DAG learning. In Proceedings of the 8th International Conference on Learning Representations.
- Maddison, C. J., Mnih, A., & Teh, Y. W. (2017). The concrete distribution: A continuous relaxation of discrete random variables. Proceedings of the 34th International Conference on Machine Learning.
- Ng, I., Fang, Z., Zhu, S., Chen, Z., & Wang, J. (2019). Masked gradient-based causal structure learning. arXiv preprint arXiv:1910.08527.
- Pamfil, R., Sriwattanaworachai, N., Desai, S., Pilgerstorfer, P., Georgatzis, K., Beaumont, P., & Aragam, B. (2020). Dynotears: Structure learning from time-series data. In International Conference on Artificial Intelligence and Statistics (pp. 1595–1605).: PMLR.
- Pearl, J. (2009). Causality: Models, Reasoning and Inference. Cambridge University Press, 2nd edition.
- Peters, J., Janzing, D., & Schölkopf, B. (2017). Elements of Causal Inference Foundations and Learning Algorithms. MIT Press.
- Peters, J., M. Mooij, J., Janzing, D., & Schölkopf, B. (2014). Causal discovery with continuous additive noise models. Journal of Machine Learning Research.
- Shpitser, I. & Pearl, J. (2006). Identification of joint interventional distributions in recursive semi-markovian causal models. In Proceedings of the National Conference on Artificial Intelligence, volume 21 (pp. 1219).: Menio Park, CA; Cambridge, MA; London; AAAI Press; 1999.
- Verma, T. & Pearl, J. (1991). Equivalence and synthesis of causal models. UCLA, Computer Science Department.
- Yang, K. D., Katcoff, A., & Uhler, C. (2018). Characterizing and learning equivalence classes of causal DAGs under interventions. Proceedings of the 35th International Conference on Machine Learning.
- Zheng, X., Aragam, B., Ravikumar, P., & Xing, E. (2018). Dags with no tears: Continuous optimization for structure learning. In Advances in Neural Information Processing Systems 31.