



Learning Causal Structures via Gradient-Based Optimization

Sébastien Lachapelle

Mila, Université de Montréal

March 4th, 2020

Overview

- Causality Framework
 - Causal Graphical Models
 - Motivating example
 - Markov Equivalence and Structure Identifiability
- Causal Structure Learning
 - Problem formulation
 - Discrete Search Algorithms
 - Gradient-Based Algorithms
- GraN-DAG & extensions
 - The algorithm
 - With interventional data
 - Neural Autoregressive Flows

Causal graphical models (CGM)

- Random vector $X \in \mathbb{R}^d$ (d variables)
- Let \mathcal{G} be a **directed acyclic graph** (DAG)
- Assume $p(x) = \prod_{i=1}^d p(x_i | x_{\pi_i^{\mathcal{G}}})$
 $\pi_i^{\mathcal{G}}$ = parents of i in \mathcal{G}
- Encodes **(conditional) independence statements**
 (via d -separation, see [Koller & Friedman, 2009])
- Almost identical to Bayesian Networks but allows
 for **interventional distributions**:
 $p(x | do(z))$

Simple example

$$\mathcal{G} = (V, E)$$



$$p(x, y, z) = p(x)p(z | x)p(y | z)$$

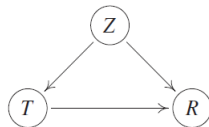
$$\implies p(x, y | z) = p(x | z)p(y | z)$$

i.e. $X \perp\!\!\!\perp Y \mid Z$

The **do operator** will be explained in the following example...

Why should you care: Kidney Stone Treatment

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



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

(Example taken from *Element of Causal Inference* by Peters et al. p111)

Why should you care: Kidney Stone Treatment

Pay attention to these two questions...
Assuming the size of your stone is unknown...

Why should you care: Kidney Stone Treatment

Pay attention to these two questions...
Assuming the size of your stone is unknown...

What is your chance of recovery knowing that the doctor gave you treatment A?

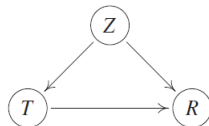
What is your chance of recovery if you decide to take treatment A?

Why should you care: Kidney Stone Treatment

$T = \text{Treatment} \in \{A, B\}$

$Z = \text{Stone size} \in \{\text{small}, \text{large}\}$

$R = \text{Patient recovered} \in \{0, 1\}$



What is your chance of recovery knowing that the doctor gave you treatment A?

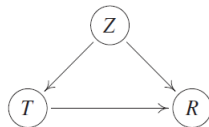
- Knowing that your doctor gave you treatment A tells you that you probably have a large kidney stone ... $P(Z = \text{large} | T = A) = 0.75$
- ... which reduces your chance of recovery
 $P(R = 1 | T = A, Z = \text{large}) = 0.73 < 0.93 = P(R = 1 | T = A, Z = \text{small})$

Why should you care: Kidney Stone Treatment

T = Treatment $\in \{A, B\}$

Z = Stone size $\in \{\text{small}, \text{large}\}$

R = Patient recovered $\in \{0, 1\}$



What is your chance of recovery knowing that the doctor gave you treatment A?

- Knowing that your doctor gave you treatment A tells you that you probably have a large kidney stone ... $P(Z = \text{large} | T = A) = 0.75$
- ... which reduces your chance of recovery
 $P(R = 1 | T = A, Z = \text{large}) = 0.73 < 0.93 = P(R = 1 | T = A, Z = \text{small})$

What is your chance of recovery if you decide to take treatment A?

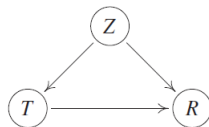
- You really don't know anything about your kidney stone
- **You taking treatment A is not a function of any variable**

Why should you care: Kidney Stone Treatment

$T = \text{Treatment} \in \{A, B\}$

$Z = \text{Stone size} \in \{\text{small, large}\}$

$R = \text{Patient recovered} \in \{0, 1\}$



What is your chance of recovery knowing that the doctor gave you treatment A?

$$P(R = 1 | T = A) = 0,78$$

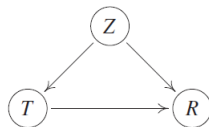
$$P(R = 1 | T = B) = \mathbf{0,83}$$

Why should you care: Kidney Stone Treatment

T = Treatment $\in \{A, B\}$

Z = Stone size $\in \{\text{small, large}\}$

R = Patient recovered $\in \{0, 1\}$



What is your chance of recovery knowing that the doctor gave you treatment A?

$$P(R = 1 | T = A) = 0,78$$

$$P(R = 1 | T = B) = \mathbf{0,83}$$

What is your chance of recovery if you decide to take treatment A?

$$P(R = 1 | do(T = A)) = \mathbf{0,832}$$

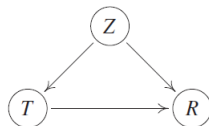
$$P(R = 1 | do(T = B)) = 0,782$$

Why should you care: Kidney Stone Treatment

T = Treatment $\in \{A, B\}$

Z = Stone size $\in \{\text{small, large}\}$

R = Patient recovered $\in \{0, 1\}$



What is your chance of recovery knowing that the doctor gave you treatment A?

$$P(R = 1 | T = A) = 0,78$$

$$P(R = 1 | T = B) = \mathbf{0,83}$$

What is your chance of recovery if you decide to take treatment A?

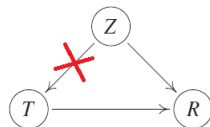
$$P(R = 1 | do(T = A)) = \mathbf{0,832}$$

$$P(R = 1 | do(T = B)) = 0,782$$

But how do we compute these **interventional distributions**?!

Why should you care: Kidney Stone Treatment

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



$$P(R, Z | do(T = A)) = P(R | Z, T = A) \underbrace{P(T = A | Z)} P(Z)$$

The decision of taking treatment A does not depend on Z anymore

Then simply marginalize as usual:

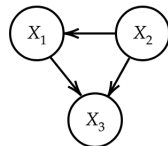
$$\begin{aligned}
 P(R = 1 | do(T = A)) &= \sum_Z P(R = 1, Z | do(T = A)) \\
 &= \sum_Z P(R = 1 | Z, T = A) P(Z) = 0,832
 \end{aligned}$$

Structure Learning

- In the kidney stone example, **the causal graph was known**
- What if we don't have it? **Learn it!**

Purely observational data

	X_1	X_2	X_3
sample 1	1.76	10.46	0.002
sample2	3.42	78.6	0.011
...
sample n	4.56	9.35	1.96

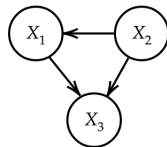
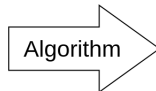


Structure Learning

- In the kidney stone example, **the causal graph was known**
- What if we don't have it? **Learn it!**

Purely observational data

	X_1	X_2	X_3
sample 1	1.76	10.46	0.002
sample 2	3.42	78.6	0.011
...
sample n	4.56	9.35	1.96



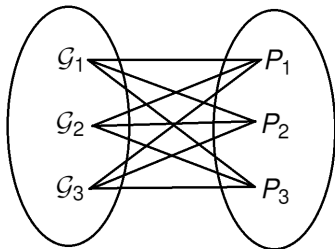
Is it even possible?

Identifiability

- In general, this is impossible without interventional data...

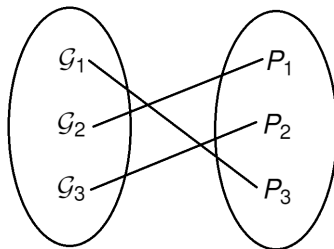
Identifiability

- In general, this is impossible without interventional data...
- Multiple DAGs can express the same distribution...



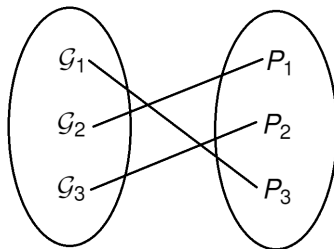
Identifiability

- If we assume causal mechanisms are "simple", then \mathcal{G} can be identified...



Identifiability

- If we assume causal mechanisms are "simple", then \mathcal{G} can be identified...



An example (useful later!)

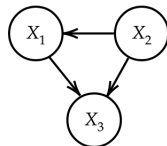
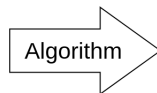
If data follows this model...

$$X_i | X_{\pi_i^{\mathcal{G}}} \sim \mathcal{N}(f_i(X_{\pi_i^{\mathcal{G}}}), \sigma_i^2)$$

...then **correct causal DAG \mathcal{G} can be identified from purely observational data** (see [Peters et al., 2014] for proof and regularity conditions)

Structure Learning

	X_1	X_2	X_3
sample 1	1.76	10.46	0.002
sample 2	3.42	78.6	0.011
...		...	
sample n	4.56	9.35	1.96



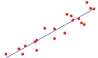
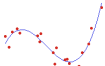
Score-based algorithms

$$\hat{\mathcal{G}} = \arg \max_{\mathcal{G} \in \text{DAG}} \text{Score}(\mathcal{G})$$

Often, $\text{Score}(\mathcal{G}) =$ regularized maximum likelihood under \mathcal{G}

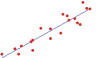
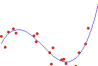
Structure Learning

Taxonomy of score-based algorithms (non-exhaustive)

		Discrete optim.	Continuous optim.
	Linear	GES [Chickering, 2003]	NOTEARS [Zheng et al., 2018]
	Nonlinear	CAM [Bühlmann et al., 2014]	GraN-DAG [Lachapelle et al., 2020]

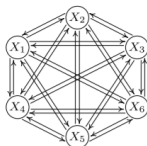
Structure Learning

Taxonomy of score-based algorithms (non-exhaustive)

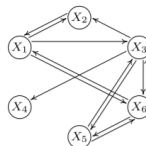
		Discrete optim.	Continuous optim.
	Linear	GES [Chickering, 2003]	NOTEARS [Zheng et al., 2018]
	Nonlinear	CAM [Bühlmann et al., 2014]	GraN-DAG [Lachapelle et al., 2020]

A greedy algorithm - CAM [Bühlmann et al., 2014]

Step 1
Choose potential parents



PNS

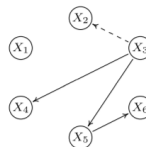


Step 2
- Start with empty graph
- Add edge improving likelihood the most
- Repeat

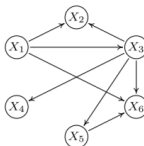
Score improvements matrix

-	0.2	0.1	-	-	0.3
0.4	-	-	-	-	-
-	0.6	-	-	-	0.4
-	-	-	-	-	-
0.3	-	-	-	-	-

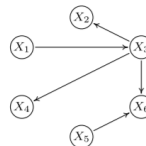
include best edge



Step 3
Remove edges via feature selection (significance test)



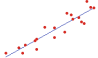
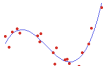
remove edges
by variable selection



Figures from [Bühlmann et al., 2014]

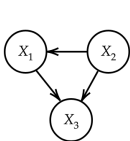
Structure Learning

Taxonomy of score-based algorithms (non-exhaustive)

		Discrete optim.	Continuous optim.
	Linear	GES [Chickering, 2003]	NOTEARS [Zheng et al., 2018]
	Nonlinear	CAM [Bühlmann et al., 2014]	GraN-DAG [Lachapelle et al., 2020]

NOTEARS: Continuous optimization for structure learning

- Encode graph as a **weighted adjacency matrix** $U = [u_1 | \dots | u_d] \in \mathbb{R}^{d \times d}$



$$A = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$$

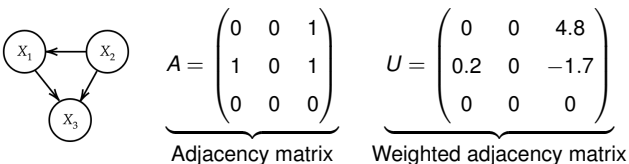
Adjacency matrix

$$U = \begin{pmatrix} 0 & 0 & 4.8 \\ 0.2 & 0 & -1.7 \\ 0 & 0 & 0 \end{pmatrix}$$

Weighted adjacency matrix

NOTEARS: Continuous optimization for structure learning

- Encode graph as a **weighted adjacency matrix** $U = [u_1 | \dots | u_d] \in \mathbb{R}^{d \times d}$

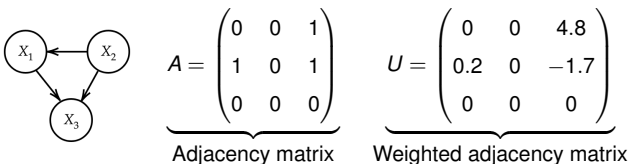


- Represents coefficients in a **linear model**:

$$X_i := u_i^\top X + \text{noise}_i \quad \forall i$$

NOTEARS: Continuous optimization for structure learning

- Encode graph as a **weighted adjacency matrix** $U = [u_1 | \dots | u_d] \in \mathbb{R}^{d \times d}$



- Represents coefficients in a **linear model**:

$$X_i := u_i^\top X + \text{noise}_i \quad \forall i$$

- For an arbitrary U , associated graph might be cyclic

Acyclicity constraint

NOTEARS [Zheng et al., 2018] uses this **differentiable acyclicity constraint**:

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

NOTEARS: Continuous optimization for structure learning

- NOTEARS [Zheng et al., 2018]:

Solve this **continuous constrained optimization problem**:

$$\max_U \underbrace{-\|\mathbf{X} - \mathbf{X}U\|_F^2 - \lambda\|U\|_1}_{\text{Score}} \quad \text{s.t.} \quad \text{Tr} e^{U \odot U} - d = 0$$

- where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the design matrix containing all n samples

NOTEARS: Continuous optimization for structure learning

- NOTEARS [Zheng et al., 2018]:

Solve this **continuous constrained optimization problem**:

$$\max_U \underbrace{-\|\mathbf{X} - \mathbf{X}U\|_F^2 - \lambda\|U\|_1}_{\text{Score}} \quad \text{s.t.} \quad \text{Tr } e^{U \odot U} - d = 0$$

- where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the design matrix containing all n samples
- Solve approximately using an **Augmented Lagrangian method**
- Amounts to maximizing (with gradient ascent)

$$-\|\mathbf{X} - \mathbf{X}U\|_F^2 - \lambda\|U\|_1 - \alpha_t(\text{Tr } e^{U \odot U} - d) - \frac{\mu_t}{2}(\text{Tr } e^{U \odot U} - d)^2$$

- while gradually increasing α_t and μ_t

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

$(A^k)_{ii}$ = number of **cycles** of length k passing through i

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

$(A^k)_{ii}$ = number of **cycles** of length k passing through i

Graph acyclic $\iff (A^k)_{ii} = 0$ for all i and all k

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

$(A^k)_{ii}$ = number of **cycles** of length k passing through i

Graph acyclic $\iff (A^k)_{ii} = 0$ for all i and all k

$$\iff \text{Tr} \left[\sum_{k=1}^{\infty} \frac{A^k}{k!} \right] = 0$$

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

$(A^k)_{ii}$ = number of **cycles** of length k passing through i

Graph acyclic $\iff (A^k)_{ii} = 0$ for all i and all k

$$\iff \text{Tr} \left[\sum_{k=1}^{\infty} \frac{A^k}{k!} \right] = 0$$

$$\iff \text{Tr} \left[\sum_{k=0}^{\infty} \frac{A^k}{k!} - A^0 \right] = 0$$

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

$(A^k)_{ii}$ = number of **cycles** of length k passing through i

Graph acyclic $\iff (A^k)_{ii} = 0$ for all i and all k

$$\iff \text{Tr} \left[\sum_{k=1}^{\infty} \frac{A^k}{k!} \right] = 0$$

$$\iff \text{Tr} \left[\sum_{k=0}^{\infty} \frac{A^k}{k!} - A^0 \right] = 0$$

$$\iff \text{Tr } e^A - d = 0$$

NOTEARS: The acyclicity constraint

$$\text{Tr } e^{U \odot U} - d = 0 \quad \left(e^M \triangleq \sum_{k=0}^{\infty} \frac{M^k}{k!} \right)$$

Suppose $A \in \{0, 1\}^{d \times d}$ is an adjacency matrix for a certain directed graph

$(A^k)_{ii}$ = number of **cycles** of length k passing through i

Graph acyclic $\iff (A^k)_{ii} = 0$ for all i and all k

$$\iff \text{Tr} \left[\sum_{k=1}^{\infty} \frac{A^k}{k!} \right] = 0$$

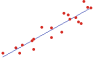
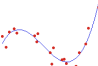
$$\iff \text{Tr} \left[\sum_{k=0}^{\infty} \frac{A^k}{k!} - A^0 \right] = 0$$

$$\iff \text{Tr } e^A - d = 0$$

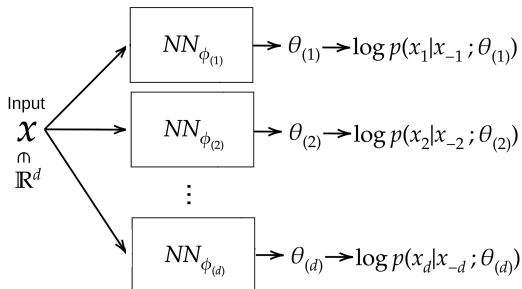
The argument is almost identical when using weighted adjacency U instead of A ...

Structure Learning

Taxonomy of score-based algorithms (non-exhaustive)

		Discrete optim.	Continuous optim.
	Linear	GES [Chickering, 2003]	NOTEARS [Zheng et al., 2018]
	Nonlinear	CAM [Bühlmann et al., 2014]	GraN-DAG [Lachapelle et al., 2020]

Gradient-Based Neural DAG Learning

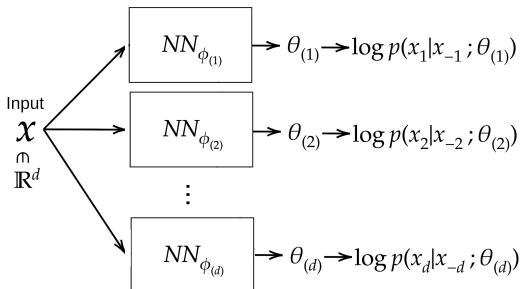


$$\phi(i) \triangleq \{W_{(i)}^{(1)}, \dots, W_{(i)}^{(L+1)}\}$$

$$W_{(i)}^{(\ell)} = \ell\text{th weight matrix of } NN_{\phi(i)}$$

$$\phi \triangleq \{\phi(i)\}_{i=1}^d$$

Gradient-Based Neural DAG Learning



$$\phi_{(i)} \triangleq \{W_{(i)}^{(1)}, \dots, W_{(i)}^{(L+1)}\}$$

$W_{(i)}^{(\ell)}$ = ℓ th weight matrix of $NN_{\phi_{(i)}}$

$$\phi \triangleq \{\phi_{(i)}\}_{i=1}^d$$

$\prod_{i=1}^d p(x_i | x_{-i}; \theta_{(i)})$ does not decompose according to a DAG!

We need to constrain the networks to be acyclic! How?

Gradient-Based Neural DAG Learning

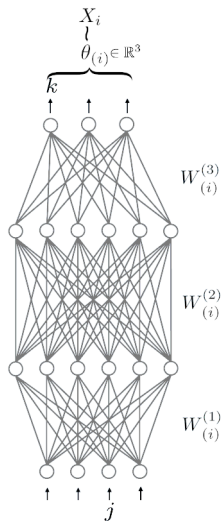
Key idea:

Construct a **weighted adjacency matrix** A_ϕ (analogous to U from the linear case) which could be used in the acyclicity constraint

Then maximize likelihood under acyclicity constraint via **augmented Lagrangian**

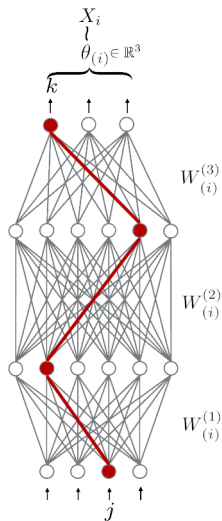
$$\max_{\phi} \underbrace{\mathbb{E}_{X \sim P_X} \sum_{i=0}^d \log p_{\phi}(X_i | X_{-i}) - \alpha_t (\text{Tr } e^{A_\phi} - d) - \frac{\mu t}{2} (\text{Tr } e^{A_\phi} - d)^2}_{\text{Augmented Lagrangian}}$$

Constructing weighted adjacency matrix A_ϕ



Let's measure the "strength" of edge $X_j \rightarrow X_i$

Constructing weighted adjacency matrix A_ϕ

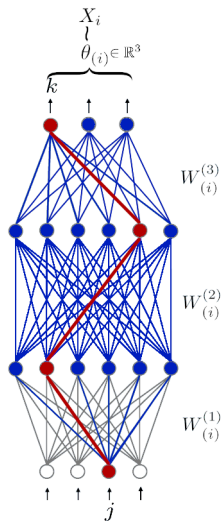


Let's measure the "strength" of edge $X_j \rightarrow X_i$

■ **Path product:**

$$|W_{h_1 j}^{(1)}| |W_{h_2 h_1}^{(2)}| |W_{k h_2}^{(3)}| \geq 0$$

Constructing weighted adjacency matrix A_ϕ



Let's measure the "strength" of edge $X_j \rightarrow X_i$

- **Path product:**

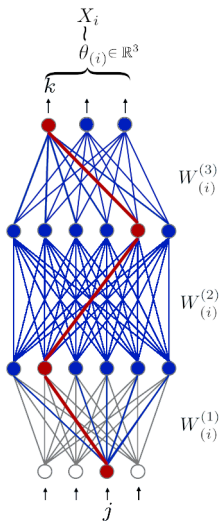
$$|W_{h_1 j}^{(1)}| |W_{h_2 h_1}^{(2)}| |W_{k h_2}^{(3)}| \geq 0$$

- $C \triangleq |W^{(3)}| |W^{(2)}| |W^{(1)}|$

"Connection strength" from X_j to $\theta_{(i)}$:

$$\sum_{k=1}^m C_{kj} \geq 0$$

Constructing weighted adjacency matrix A_ϕ



Let's measure the "strength" of edge $X_j \rightarrow X_i$

- **Path product:**

$$|W_{h_1 j}^{(1)}| |W_{h_2 h_1}^{(2)}| |W_{k h_2}^{(3)}| \geq 0$$

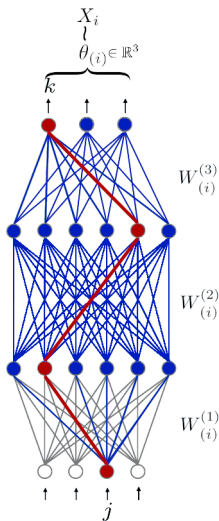
- $C \triangleq |W^{(3)}| |W^{(2)}| |W^{(1)}|$

"Connection strength" from X_j to $\theta_{(i)}$:

$$\sum_{k=1}^m C_{kj} \geq 0$$

- $\sum_{k=1}^m C_{kj} = 0 \Rightarrow$ All paths from X_j to X_i are **inactive!**

Constructing weighted adjacency matrix A_ϕ



Let's measure the "strength" of edge $X_j \rightarrow X_i$

- **Path product:**

$$|W_{h_1 j}^{(1)}| |W_{h_2 h_1}^{(2)}| |W_{k h_2}^{(3)}| \geq 0$$

- $C \triangleq |W^{(3)}| |W^{(2)}| |W^{(1)}|$

"Connection strength" from X_j to $\theta_{(i)}$:

$$\sum_{k=1}^m C_{kj} \geq 0$$

- $\sum_{k=1}^m C_{kj} = 0 \Rightarrow$ All paths from X_j to X_i are **inactive!**

$$(A_\phi)_{ji} \triangleq \begin{cases} \sum_{k=1}^m (C^{(i)})_{kj}, & \text{if } i \neq j \\ 0, & \text{otherwise} \end{cases}$$

Gradient-Based Neural DAG Learning

The algorithm:

- 1 Preliminary neighborhood selection (analogous to CAM)
i.e. for each node, select potential parents via any *variable selection* approach
- 2 Maximize likelihood under acyclicity constraint via **augmented Lagrangian**

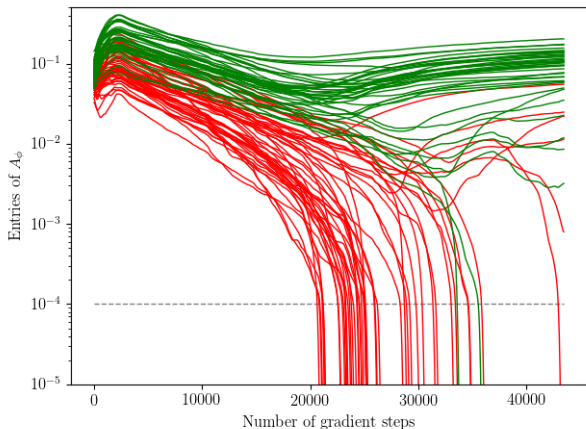
$$\max_{\phi} \mathbb{E}_{X \sim P_X} \sum_{i=0}^d \log p_{\phi}(x_i | x_{-i}) - \alpha_t (\text{Tr } e^{A_{\phi}} - d) - \frac{\mu_t}{2} (\text{Tr } e^{A_{\phi}} - d)^2$$

Augmented Lagrangian

- 3 DAG Pruning (analogous to CAM)
i.e. for each node, get rid of some parents via any *variable selection* approach

Step 1 and 3 helps reducing overfitting.
Important since **adding edges cannot reduce maximum likelihood**

Gradient-Based Neural DAG Learning



Correct edges
Wrong edges

Experiments

Synthetic data: $X_i | X_{\pi_i^G} \sim \mathcal{N}(f_i(X_{\pi_i^G}), \sigma_i^2)$ $f_i \sim$ Gaussian Process

Models: GraN-DAG, NOTEARS and CAM makes the Gaussian assumption

Real data: Measurements of expression levels of proteins and phospholipids in human immune system cells [Sachs et al., 2005]

		Synthetic (50 nodes)		Protein data set	
		SHD	SID	SHD	SID
Continuous	GraN-DAG	102.6±21.2	1060.1±109.4	13	47
	DAG-GNN	191.9±15.2	2146.2±64	16	44
	NOTEARS	202.3±14.3	2149.1±76.3	21	44
Discrete	CAM	98.8±20.7	1197.2±125.9	12	55
	RANDOM	708.4±234.4	1921.3±203.5	21	60

DAG-GNN [Yu et al., 2019]

Experiments

- In previous setup, synthetic data generation and model matched
- Here: **model misspecification**

Synthetic post nonlinear data sets

		PNL-GP		PNL-MULT	
		SHD	SID	SHD	SID
10 nodes ER1	GraN-DAG	1.6±3.0	3.9±8.0	13.1±3.8	35.7±12.3
	DAG-GNN	11.5±6.8	32.4±19.3	17.900±6.2	40.700±14.743
	NOTEARS	10.7±5.5	34.4±19.1	14.0±4.0	38.6±11.9
	CAM	1.5±2.6	6.8±12.1	12.0±6.4	36.3±17.7
	GSF	6.2±3.3	[7.7±8.7, 18.9±12.4]	10.7±3.0	[9.8±11.9, 25.3±11.5]
	RANDOM	23.8±2.9	36.8±19.1	23.7±2.9	37.7±20.7

GSF [Huang et al., 2018a]

Experiments: Effect of sample size

- Previous experiment: relatively small dataset: 1000 examples
- GraN-DAG is more expressive than CAM
- Advantage shows up in large sample size regimes

Effect of sample size - Gauss-ANM 50 nodes ER4 (averaged over 10 datasets)

Sample size	Method	SHD	SID
500	CAM	123.5 ± 13.9	1181.2 ± 160.8
	GraN-DAG	130.2 ± 14.4	1246.4 ± 126.1
1000	CAM	103.7 ± 15.2	1074.7 ± 125.8
	GraN-DAG	104.4 ± 15.3	942.1 ± 69.8
5000	CAM	74.1 ± 13.2	845.0 ± 159.8
	GraN-DAG	71.9 ± 15.9	554.1 ± 117.9
10000	CAM	66.3 ± 16.0	808.1 ± 142.9
	GraN-DAG	65.9 ± 19.8	453.4 ± 171.7

GraN-DAG with interventions [Brouillard et al., 2020]

Can we make use of **interventional data**?

GraN-DAG with interventions [Brouillard et al., 2020]

Some **terminology** and **setting**:

- $I \subset \{1, \dots, n\}$ is an *interventional target* (set of nodes on which we intervene)
- **Definition** of *stochastic intervention*:

$$p(x_1, \dots, x_d | do(X_I)) \triangleq \prod_{j \notin I} p_j(x_j | x_{\pi_j^G}) \prod_{j \in I} \tilde{p}_j(x_j)$$

where $\tilde{p}_j(x_j)$ is the *new marginal* replacing $p_j(x_j | x_{\pi_j^G})$ (parents are "cut out")

GraN-DAG with interventions [Brouillard et al., 2020]

Some **terminology** and **setting**:

- $I \subset \{1, \dots, n\}$ is an *interventional target* (set of nodes on which we intervene)

- **Definition** of *stochastic intervention*:

$$p(x_1, \dots, x_d | do(X_I)) \triangleq \prod_{j \notin I} p_j(x_j | x_{\pi_j^G}) \prod_{j \in I} \tilde{p}_j(x_j)$$

where $\tilde{p}_j(x_j)$ is the *new marginal* replacing $p_j(x_j | x_{\pi_j^G})$ (parents are "cut out")

- **Observed:** $\{(X^{(1)}, I^{(1)}), \dots, (X^{(n)}, I^{(n)})\}$ where $I^{(i)}$ is the interventional target associated to observation $X^{(i)}$.

$$\begin{aligned} I^{(i)} &\sim P(I) \text{ i.i.d. } \forall i \\ X^{(i)} | I^{(i)} &\sim P(X | I = I^{(i)}) \triangleq p(x_1, \dots, x_d | do(X_{I^{(i)}})) \forall i \end{aligned} \quad (1)$$

where $P(I)$ is a distribution over a collection of interventional targets \mathcal{I}

GraN-DAG with interventions [Brouillard et al., 2020]

- Think about a CGM as a family of models of the form

$$\left\{ \prod_{j \notin I} p_j(x_j | x_{\pi_j^G}; \phi_j) \prod_{j \in I} \tilde{p}_j(x_j; \omega_j^I) \mid I \in \mathcal{I} \right\}$$

where $\omega^I \triangleq \{\omega_j^I\}_{j \in I}$ for each $I \in \mathcal{I}$ are learnable parameters.

- The natural **optimization problem**:

$$\max_{\phi, \{\omega^I\}_{I \in \mathcal{I}}} \mathbb{E}_{(X, I) \sim P(X, I)} \left[\sum_{j \notin I} \log p_j(X_j | X_{-j}; \phi_j) + \sum_{j \in I} \log p_j(X_j; \omega_j^I) \right] \quad \text{s.t.} \quad \text{Tr } e^{A_\phi} = d$$

- But we do not really care about learning the $p_j(X_j; \omega_j^I)$...
- ... and problem trivially decomposes as a sum of \max_{ϕ} and $\max_{\{\omega^I\}_{I \in \mathcal{I}}}$ so ...

GraN-DAG with interventions [Brouillard et al., 2020]

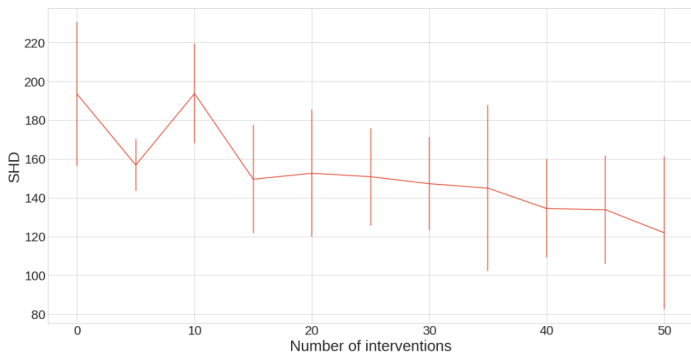
- ... can forget about $p_j(X_j; \omega_j^I)$ altogether and get

The optimization problem:

$$\max_{\phi} \mathbb{E}_{(X,I) \sim P(X,I)} \sum_{j \notin I} \log p(X_j | X_{-j}; \phi_j) \quad \text{s.t.} \quad \text{Tr } e^{A_{\phi}} = d$$

In a nutshell: We throw out the conditionals associated with the intervention variables

GraN-DAG with interventions [Brouillard et al., 2020]



- Linear data (unidentifiable without interventions)
- 50 nodes and ≈ 200 edges
- Intervention on one node at a time

GraN-DAG with interventions [Brouillard et al., 2020]

■ Nonlinear data

Method	20 nodes, $e = 1$		20 nodes, $e = 4$		50 nodes, $e = 1$		50 nodes, $e = 4$	
	SHD	SID	SHD	SID	SHD	SID	SHD	SID
GraN-DAG	1.0 \pm 1.2	1.9 \pm 3.7	33.3 \pm 10.0	138.9 \pm 21.0	85.3 \pm 20.2	885.3 \pm 151.4	3.2 \pm 2.8	12.7 \pm 13.6
GraN-DAG no interv	0.7 \pm 0.8	1.1 \pm 2.3	41.5 \pm 8.0	164.2 \pm 27.9	109.8 \pm 17.6	1021.7 \pm 109.0	4.1 \pm 2.3	18.6 \pm 15.1
GIES	14.6 \pm 5.8	34.9 \pm 23.1	64.3 \pm 6.1	282.7 \pm 36.7	170.2 \pm 21.1	1820.7 \pm 183.0	41.6 \pm 9.9	107.9 \pm 48.9
CAM*	1.9 \pm 2.5	4.8 \pm 6.3	48.8 \pm 23.5	144.3 \pm 38.6	91.9 \pm 11.9	1024.4 \pm 118.0	4.7 \pm 3.8	24.5 \pm 18.1

■ Linear data

Method	20 nodes, $e = 1$		20 nodes, $e = 4$		50 nodes, $e = 1$		50 nodes, $e = 4$	
	SHD	SID	SHD	SID	SHD	SID	SHD	SID
GraN-DAG	5.7 \pm 5.4	31.3 \pm 37.5	24.5 \pm 7.4	159.6 \pm 34.9	12.5 \pm 7.2	63.4 \pm 36.0	52.7 \pm 16.0	699.9 \pm 166.3
GraN-DAG no interv	17.7 \pm 8.2	91.8 \pm 71.9	88.3 \pm 23.2	275.2 \pm 17.0	41.2 \pm 7.3	163.8 \pm 70.4	193.4 \pm 39.3	1667.7 \pm 169.2
GIES	2.4 \pm 1.1	0.0 \pm 0.0	29.7 \pm 28.2	122.3 \pm 103.7	13.9 \pm 3.8	0.0 \pm 0.0	129.7 \pm 75.1	757.8 \pm 477.7
CAM*	6.5 \pm 8.7	19.0 \pm 35.0	59.2 \pm 27.2	173.1 \pm 69.2	1.9 \pm 2.4	9.8 \pm 16.5	135.4 \pm 29.8	1484.1 \pm 274.8

■ More experiments in workshop paper...

GraN-DAG with Neural Autoregressive flows

- In previous experiments, GraN-DAG models was:

$$X_i = NN_{\phi_i}(X_{\pi_i^G}) + \sigma_i Z \quad \text{with} \quad Z \sim \mathcal{N}(0, 1) \quad \forall i$$

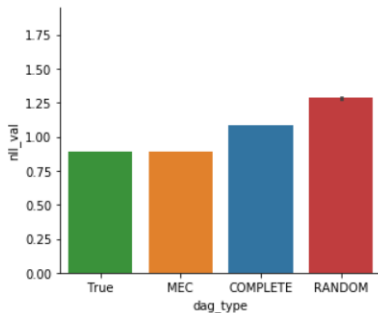
- GraN-DAG's framework allows for usage of "Neural Autoregressive Flows" [Huang et al., 2018b]

$$X_i = NAF(Z; NN_{\phi_i}(X_{\pi_i^G})) \quad \text{with} \quad Z \sim \mathcal{N}(0, 1) \quad \forall i$$

- The function $NAF(\cdot; NN_{\phi_i}(X_{\pi_i^G}))$ is **invertible** and with **tractable Jacobian** so the likelihood of X can be computed exactly and maximized

GraN-DAG with Neural Autoregressive flows

- Without interventions, we run into identifiability problems ...



- Future work:** make it works with interventional data (since identifiability is less of a problem)

Conclusion and future work

Gradient-based DAG search...

- ... performs similarly to its discrete analogs
- ... scales well with number of samples (since amenable to stochastic optimization)
- ... can be easily adapted to work with interventional data
- ... allows for very expressive density models (Neural Autoregressive flow)

Future work:

- DAGs appear in many places, could we adapt the neural acyclicity constraint to other problems? (Not causality?)
- Drawing links between causality and representation learning

References

Brouillard, P., Drouin, A., Lachapelle, S., Lacoste, A., & Lacoste-Julien, S. (2020). Gradient-based neural dag learning with interventions.

Bühlmann, P., Peters, J., & Ernest, J. (2014). CAM: Causal additive models, high-dimensional order search and penalized regression. *Annals of Statistics*.

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

Huang, B., Zhang, K., Lin, Y., Schölkopf, B., & Glymour, C. (2018a). Generalized score functions for causal discovery. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*.

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

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 Eighth International Conference on Learning Representations (to appear)*.

Peters, J., M. Mooij, J., Janzing, D., & Schölkopf, B. (2014). Causal discovery with continuous additive noise models. *Journal of Machine Learning Research*.

Sachs, K., Perez, O., Pe'er, D., Lauffenburger, D., & Nolan, G. (2005). Causal protein-signaling networks derived from multiparameter single-cell data. *Science*.

Yu, Y., Chen, J., Gao, T., & Yu, M. (2019). DAG-GNN: DAG structure learning with graph neural networks. In *Proceedings of the 36th 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*.



Philippe Brouillard



Alexandre Drouin



Tristan Deleu



Alexandre Lacoste



Simon Lacoste-Julien

