

# Causal structure learning and sampling using Markov Monte Carlo with momentum

---

MORITZ SCHAUER

*Chalmers University of Technology | University of Gothenburg*

With MARCEL WIENÖBST, UNIVERSITY OF LÜBECK.

Cramér Society 2023

# **Alzheimer's Disease Neuroimaging Initiative**

---

# Alzheimer's Disease Neuroimaging Initiative (ADNI)



1

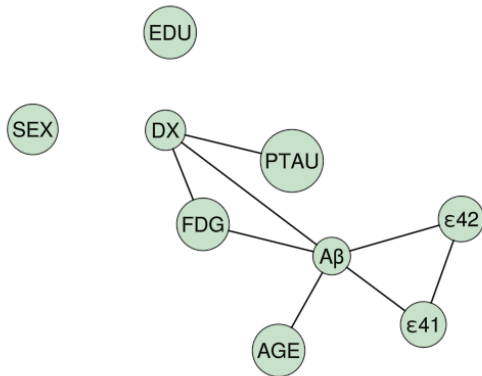
ADNI is a longitudinal study designed to develop clinical, imaging, genetic, and biochemical biomarkers for the early detection and tracking of Alzheimer's disease (AD).

<sup>1</sup>Image: ADNI.

## ADNI data

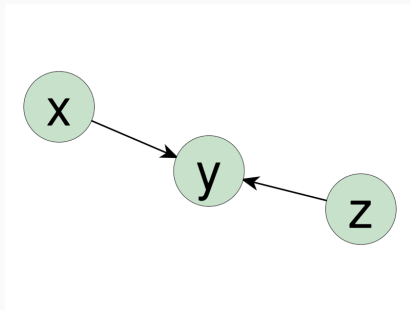
SEX	FDG	DX	$A\beta$	$\epsilon 4$	PTAU	EDU	AGE
Male	1.13615	CN	×	0	0	16	78.3
Male	1.3086	Dementia	721.5	2	22.83	18	81.3
Male	×	MCI	1501	0	13.29	10	67.5
Male	1.25956	CN	547.3	0	31.43	16	70.7
Female	×	MCI	×	0	×	13	81.4
...							

## Associations between variables



# Graphical models

---



A **DAG** is a directed graph such that following arrows it is impossible to return to any vertex (no cycles). A **PDAG** has additional undirected edges.

Vertices  $x$ ,  $y$ ,  $z$  correspond to stochastic variables.

# Classical (Bayesian) statistics

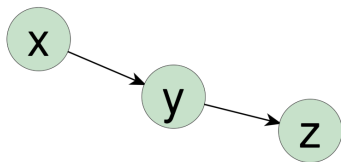
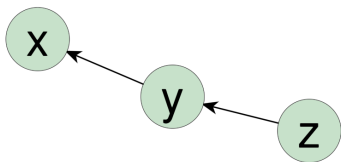
---

A single joint density on unknowns and observables. Different possible factorizations:

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



## Classical (Bayesian) statistics



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

Typically some Markovian properties, for example here

$$x \perp\!\!\!\perp z \mid y$$

and corresponding factorizations of that density.

## Faithfulness assumption

---

Assume **perfect correspondence** between law  $p$  and DAG  $\mathcal{G}$

$$x \perp\!\!\!\perp_p z \mid y \quad \Leftrightarrow \quad x \perp\!\!\!\perp_{\mathcal{G}} z \mid y$$

## Faithfulness assumption

---

Assume **perfect correspondence** between law  $p$  and DAG  $\mathcal{G}$

$$x \perp\!\!\!\perp_p z \mid y \quad \Leftrightarrow \quad x \perp\!\!\!\perp_{\mathcal{G}} z \mid y$$

For example  $\mathcal{G}$  having no edges is equivalent to complete independence under  $p$ .

## Faithfulness assumption

Assume **perfect correspondence** between law  $p$  and DAG  $\mathcal{G}$

$$x \perp\!\!\!\perp_p z \mid y \quad \Leftrightarrow \quad x \perp\!\!\!\perp_{\mathcal{G}} z \mid y$$

For example  $\mathcal{G}$  having no edges is equivalent to complete independence under  $p$ .

The implication “ $\Rightarrow$ ” is **faithfulness**.

## Faithfulness violation

$z_1, \dots, z_4$  independent noise and

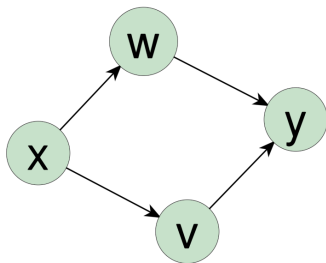
$$x = z_1$$

$$w = x + z_2$$

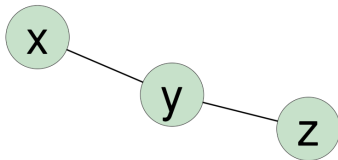
$$v = -x + z_3$$

$$y = v + w + z_4.$$

Independence  $x \perp\!\!\!\perp y$  not implied by the DAG.

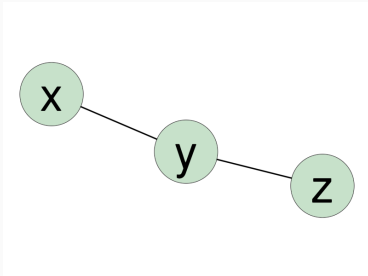


## DAG discovery



What are the possible DAG models (edge orientations) under faithfulness?

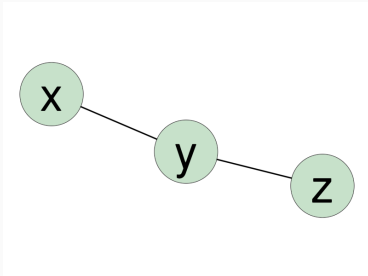
## DAG discovery



What are the possible DAG models (edge orientations) under faithfulness?

- $x \perp\!\!\!\perp z$  implies  $x \rightarrow y \leftarrow z$

## DAG discovery

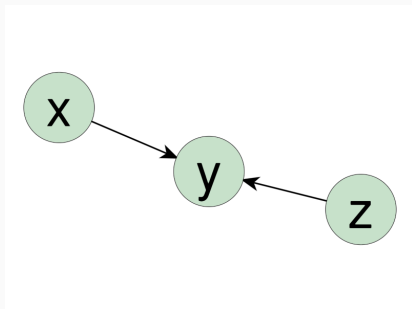


What are the possible DAG models (edge orientations) under faithfulness?

- $x \perp\!\!\!\perp z$  implies  $x \rightarrow y \leftarrow z$
- $x \not\perp\!\!\!\perp z$  is compatible with all others

$$\underbrace{x \leftarrow y \leftarrow z \quad x \rightarrow y \rightarrow z \quad x \leftarrow y \rightarrow z}_{\text{Markov equivalence class (MEC)}}$$





$x$ ,  $y$ , and  $z$  in a DAG form an **v-structure** if  $x \rightarrow y \leftarrow z$  and  $x$  and  $z$  are not adjacent.

$$p(x, y, z) = p(y \mid x, z)p(x)p(z)$$

Among babies of low birth weight ( $y$ ) maternal smoking ( $x$ ) was associated with lower infant mortality.

Among babies of low birth weight ( $y$ ) maternal smoking ( $x$ ) was associated with lower infant mortality.

Think: if  $z$  are other independent causes of low birth weight, then

$$x \rightarrow y \leftarrow z \quad \text{and} \quad x \not\perp\!\!\!\perp z \mid y.$$

## Markov equivalence classes

---

All DAGs on a vertex set  $V$  with  $n$  vertices with the same set of v-structures and the same set of adjacencies are observationally equivalent and form the **Markov equivalence class (MEC)** denoted  $\mathcal{M}_n$  (Verma and Pearl, 1990)

A way to represent a MEC is the CPDAG (completed PDAG):

Arrows  $x \rightarrow y$  only if all members of the equivalence class agree on the direction; undirected edges  $x - y$  otherwise.

A way to represent a MEC is the CPDAG (completed PDAG):

Arrows  $x \rightarrow y$  only if all members of the equivalence class agree on the direction; undirected edges  $x - y$  otherwise.

$\mathcal{M}_n$  is the space of CPDAGs or MECs with elements denoted  $\gamma, \eta, \dots \in \mathcal{M}_n$ .

## Two variables

---

What are the MECs on two variables  $x$  and  $y$ ?

## Two variables

---

What are the MECs on two variables  $x$  and  $y$ ?

No v-structures, so

$$\gamma_0 = "x \perp y" = \{ "x \perp y" \}$$

$$\gamma_1 = "x - y" = \{ "x \rightarrow y", "x \leftarrow y" \}$$



## Two variables

---

What are the MECs on two variables  $x$  and  $y$ ?

No v-structures, so

$$\gamma_0 = "x \text{ --- } y" = \{ "x \text{ --- } y" \}$$

$$\gamma_1 = "x - y" = \{ "x \rightarrow y", "x \leftarrow y" \}$$

Not so nice when trying to infer causation from association.

## Two variables

---

What are the MECs on two variables  $x$  and  $y$ ?

No v-structures, so

$$\gamma_0 = "x \quad y" = \{ "x \quad y" \}$$

$$\gamma_1 = "x - y" = \{ "x \rightarrow y", "x \leftarrow y" \}$$

Not so nice when trying to infer causation from association.

But: Re-factorizing  $p(\text{data}|\theta)p(\theta)$  as  $p(\theta|\text{data})p(\text{data})$  has some applications.

# Characterisation of CPDAGs

**THEOREM 4.1** (Characterization of  $D^*$ ). *A graph  $G \equiv (V, E)$  is equal to  $D^*$  for some ADG  $D$  if and only if  $G$  satisfies the following four conditions.*

- (i)  *$G$  is a chain graph.*
- (ii) *For every chain component  $\tau$  of  $G$ ,  $G_\tau$  is chordal.*
- (iii) *The configuration  $a \rightarrow b - c$  does not occur as an induced subgraph of  $G$ .*
- (iv) *Every arrow  $a \rightarrow b \in G$  is strongly protected in  $G$ .*

S. A. Andersson, D. Madigan and M. D. Perlman, “A characterization of Markov equivalence classes for acyclic digraphs”, *Annals of Statistics* 25 (1997) 505-541.

---

## Causal Structure Learning With Momentum: Sampling Distributions Over Markov Equivalence Classes of DAGs

---

**Moritz Schauer**

Chalmers University of Technology  
and University of Gothenburg

**Marcel Wienöbst**

Institute for Theoretical Computer Science,  
University of Lübeck

### Abstract

In the context of inferring a Bayesian network structure (directed acyclic graph, DAG for short), we devise a non-reversible continuous-time Markov chain “Causal MCMC”

### 1 Introduction

A Bayesian network is a probabilistic graphical model that represents a set of random variables and their conditional (in)dependencies using a directed acyclic graph (DAG). Graph and random variables are linked



---

## Causal Structure Learning With Momentum: Sampling Distributions Over Markov Equivalence Classes of DAGs

---

**Moritz Schauer**

Chalmers University of Technology  
and University of Gothenburg

**Marcel Wienöbst**

Institute for Theoretical Computer Science,  
University of Lübeck

### Abstract

In the context of inferring a Bayesian network structure (directed acyclic graph, DAG for short), we devise a non-reversible continuous-time Markov chain “Gibbs sampler”

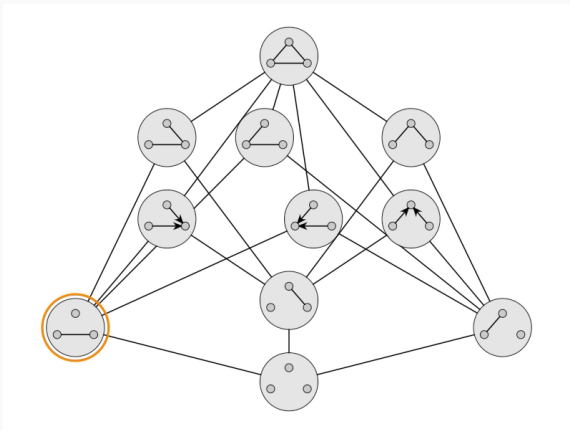
### 1 Introduction

A Bayesian network is a probabilistic graphical model that represents a set of random variables and their conditional (in)dependencies using a directed acyclic graph (DAG). Graph and random variables are linked

## Random walk on CPDAGs

---

# Random walks for causal models



## Uniform random walk on a graph

---

Random walk visits vertices with many edges more often, so needs to spend less time there:

$$\mathbb{E}[\tau_v] = \frac{1}{\deg(v)}, \quad \tau_v \text{ residence time in } v$$

After  $\tau_v$  time units, the process jumps to a neighbour (picked at random.)



## Uniform random walk on a graph

---

Random walk visits vertices with many edges more often, so needs to spend less time there:

$$\mathbb{E}[\tau_v] = \frac{1}{\deg(v)}, \quad \tau_v \text{ residence time in } v$$

After  $\tau_v$  time units, the process jumps to a neighbour (picked at random.)

With Markovianity

$$\tau_v \sim \text{Exp}(\deg(v))$$

## Neighbours??

Declare adjacency between CPDAGs:

$$\gamma = \{ "x \rightarrow y \leftarrow z", "x \leftarrow y \leftarrow z" \} (= "x - y \leftarrow z")$$

and

$$\eta = \{ x \rightarrow y \leftarrow z \}$$

are **neighbours**, because I can insert an edge into " $x \rightarrow y \leftarrow z$ " to obtain " $x \rightarrow y \leftarrow z$ ".

## Neighbours??

Declare adjacency between CPDAGs:

$$\gamma = \{ "x \rightarrow y \leftarrow z", "x \leftarrow y \leftarrow z" \} (= "x - y - z")$$

and

$$\eta = \{ x \rightarrow y \leftarrow z \}$$

are **neighbours**, because I can insert an edge into " $x \rightarrow y \leftarrow z$ " to obtain " $x \rightarrow y \leftarrow z$ ".

Notation:  $\eta \in \text{Insert}(\gamma)$ ,  $\gamma \in \text{Delete}(\eta)$ .

The operator  $\text{Insert}(\gamma, x, y, T)$  inserts the edge  $x \rightarrow y$  to the CPDAG  $\gamma$  and directs previously undirected edges  $t - y$  to  $t \rightarrow y$  for  $t \in T$ , such that vertices  $t \in T$  become “tails” of a v-structure  $t \rightarrow y \leftarrow x$ .

The operator  $\text{Insert}(\gamma, x, y, T)$  inserts the edge  $x \rightarrow y$  to the CPDAG  $\gamma$  and directs previously undirected edges  $t - y$  to  $t \rightarrow y$  for  $t \in T$ , such that vertices  $t \in T$  become “tails” of a v-structure  $t \rightarrow y \leftarrow x$ .

---

Fineprint: Here  $x$  and  $y$  are not adjacent and  $T$  are (undirected) neighbours of  $y$  that are not adjacent to  $x$ . The resulting PDAG is then completed.

Denote by  $NA_x(y)$  the (undirected) neighbours of  $y$  that are adjacent to  $x$ .

## Valid moves

---

Denote by  $\text{NA}_x(y)$  the (undirected) neighbours of  $y$  that are adjacent to  $x$ .

$\text{Insert}(\gamma, x, y, T)$  is a valid move, if and only if

- $\text{NA}_x(y)$  and the elements of  $T$  form a clique and
- any path from  $y$  to  $x$  without a directed edge pointing towards  $y$  (such a path is called semi-directed) contains a vertex in  $\text{NA}_x(y) \cup T$ .

## Valid moves

---

Denote by  $NA_x(y)$  the (undirected) neighbours of  $y$  that are adjacent to  $x$ .

$\text{Insert}(\gamma, x, y, T)$  is a valid move, if and only if

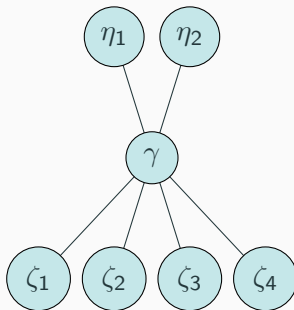
- $NA_x(y)$  and the elements of  $T$  form a clique and
- any path from  $y$  to  $x$  without a directed edge pointing towards  $y$  (such a path is called semi-directed) contains a vertex in  $NA_x(y) \cup T$ .

---

Story for the delete operator is a bit simpler

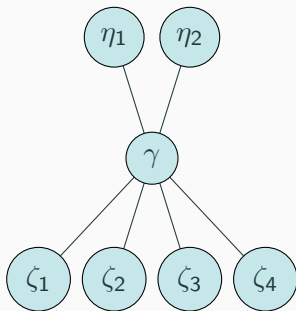


## Random walk on CPDAGs



A MEC  $\gamma$  with two neighbours  $\eta_1, \eta_2$  in  $\text{Insert}(\gamma)$  and four neighbours  $\zeta_1, \dots, \zeta_4$  in  $\text{Delete}(\gamma)$ . This is a lattice!

## Random walk on CPDAGs

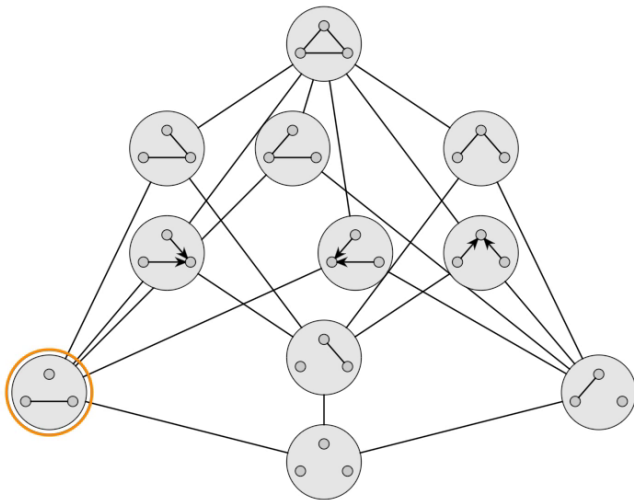


Random walk will leave  $\gamma$  after an exponentially distributed time with total rate  $\Lambda(\gamma) = 6$  towards one of the six neighbours drawn from  $\kappa_\gamma = \mathcal{U}(\{\eta_1, \eta_2, \zeta_1, \zeta_2, \zeta_3, \zeta_4\})$ . (Not so easy to count for large graphs...)

## **Adding momentum**

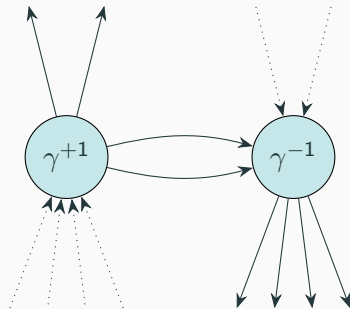
---

## Lifted random walk for causal models



## Lifted random walk

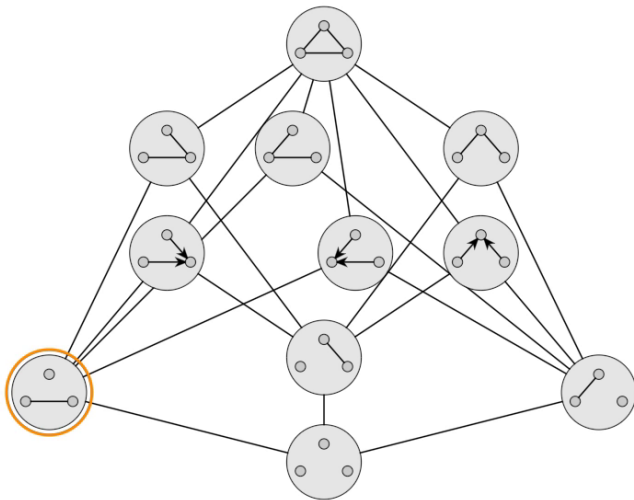
$$\gamma^{+1} := (\gamma, +1) \qquad \gamma^{-1} := (\gamma, -1)$$



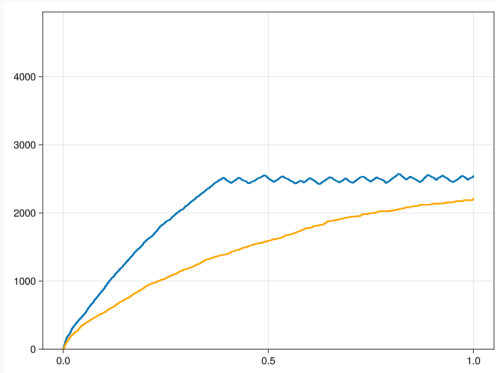
If  $\gamma \in \mathcal{M}_n$  has 2 direct neighbours in  $\text{Insert}(\gamma)$  and 4 direct neighbours in  $\text{Delete}(\gamma)$ :

Move up from  $\gamma^{+1}$  with total rate 2, move from  $\gamma^{+1}$  to  $\gamma^{-1}$  with rate  $2 = 4 - 2$  and down from  $\gamma^{-1}$  with total rate 4.

## Lifted random walk



# Mixing



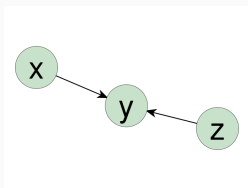
Continuous-time trace of the number of edges of the sampled graphs when targeting a uniform distribution on CPDAGs with 100 vertices. Blue: Lifted, orange: Normal.

The total time of 1 unit corresponds to 5 000 jumps.

# Causal discovery

---





One DAG and a corresponding factorization

$$p(x, y, z) = p(y \mid x, z)p(x)p(z)$$

can describe a family different of joint densities corresponding to different interventions:

$$p_{do(z=z_0)}(x, y) = p(x)p(y \mid x, z_0) \neq p(x, y)$$

$$p_{do(y=y_0)}(x, z) = p(x)p(z) = p(x, z)$$

Difficult problem: Learn a causal model from observational data.

Difficult problem: Learn a causal model from observational data.

Assuming that all relevant variables are observed, the causal model is in the observational MEC.

Difficult problem: Learn a causal model from observational data.

Assuming that all relevant variables are observed, the causal model is in the observational MEC.

If you know the MEC, you can think of experiments to pin down the causal relationships further, e.g. by gene knockouts.

## Score based causal discovery

---

## Markov equivalent score

---

A scoring function for DAGs is a **Markov equivalent score** if it assigns the same score to any DAG in the same MEC.

## Markov equivalent score

---

A scoring function for DAGs is a **Markov equivalent score** if it assigns the same score to any DAG in the same MEC.

Example: Bayesian information criterion (BIC).

## Markov equivalent score

---

A scoring function for DAGs is a **Markov equivalent score** if it assigns the same score to any DAG in the same MEC.

Example: Bayesian information criterion (BIC).

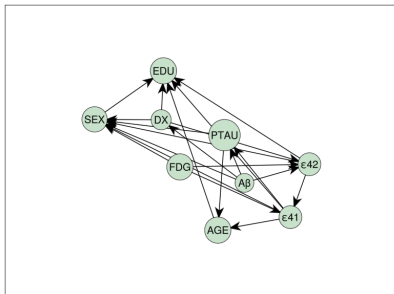
Exponentiated BIC score factorises over the DAGs

$$w(G, \text{Data}) = \prod_{x \in V} w(\text{Pa}_G(x), x, \text{Data}),$$

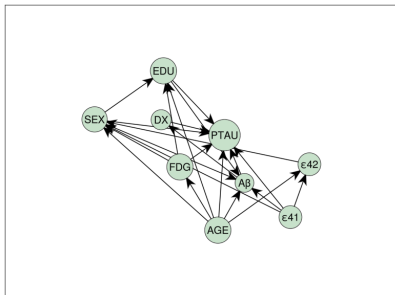
Changes in  $w$  can be computed efficiently by comparing local scores.



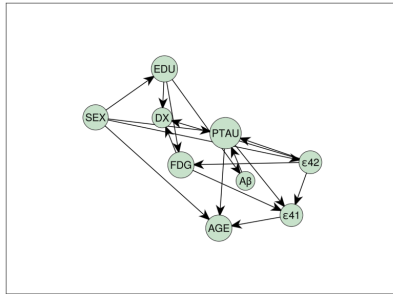
There are 1213442454842881 (1.2 quadrillion) directed acyclic graphs on 9 vertices. These are some of them.



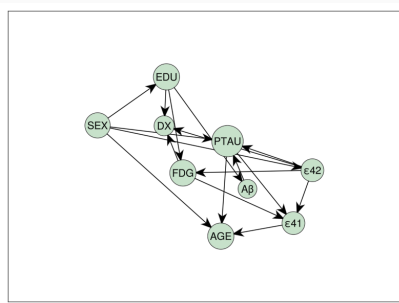
There are 1213442454842881 (1.2 quadrillion) directed acyclic graphs on 9 vertices. These are some of them.



There are 1213442454842881 (1.2 quadrillion) directed acyclic graphs on 9 vertices. These are some of them.



There are 1213442454842881 (1.2 quadrillion) directed acyclic graphs on 9 vertices. These are some of them.



Given the graph structure, perform regression on the parents of a variable to obtain a model, e.g.

$$DX = \beta_1 \cdot EDU + \beta_2 \cdot FDG + \beta_3 \cdot PTAU + \text{error term.}$$

## Zanella process

---

Continuous time random walk to sample from a distribution  $\pi$  defined on  $\mathcal{M}_n$ .

## Zanella process

---

Continuous time random walk to sample from a distribution  $\pi$  defined on  $\mathcal{M}_n$ .

Like in Metropolis-Hastings we need a balancing function  $g$  such as  $\sqrt{t}$  or  $\min(1, t)$  with the property  $g(t) = tg(1/t)$ .

## Zanella process

Continuous time random walk to sample from a distribution  $\pi$  defined on  $\mathcal{M}_n$ .

Like in Metropolis-Hastings we need a balancing function  $g$  such as  $\sqrt{t}$  or  $\min(1, t)$  with the property  $g(t) = tg(1/t)$ .

The Zanella process is defined by the jump intensity

$$\lambda(\gamma \curvearrowright \eta) = \begin{cases} g\left(\frac{\pi\{\eta\}}{\pi\{\gamma\}}\right) & \text{if } \eta \in \text{Insert}(\gamma) \sqcup \text{Delete}(\gamma) \\ 0 & \text{otherwise} \end{cases},$$

where  $\gamma \in \mathcal{M}_n$ .

## Lifted Zanella process

---

A random walk on  $\mathcal{M}_n \times \{+1, -1\}$  with correct marginal:



## Lifted Zanella process

A random walk on  $\mathcal{M}_n \times \{+1, -1\}$  with correct marginal:

For  $\gamma \in \mathcal{M}_n$ ,

$$\lambda(\gamma^{+1} \curvearrowright \eta^{+1}) = \begin{cases} g\left(\frac{\pi\{\eta\}}{\pi\{\gamma\}}\right) & \text{if } \eta \in \text{Insert}(\gamma) \\ 0 & \text{otherwise.} \end{cases}$$

## Lifted Zanella process

A random walk on  $\mathcal{M}_n \times \{+1, -1\}$  with correct marginal:

For  $\gamma \in \mathcal{M}_n$ ,

$$\lambda(\gamma^{+1} \curvearrowright \eta^{+1}) = \begin{cases} g\left(\frac{\pi\{\eta\}}{\pi\{\gamma\}}\right) & \text{if } \eta \in \text{Insert}(\gamma) \\ 0 & \text{otherwise.} \end{cases}$$

$$\lambda(\eta^{-1} \curvearrowright \gamma^{-1}) = \begin{cases} g\left(\frac{\pi\{\gamma\}}{\pi\{\eta\}}\right) & \text{if } \gamma \in \text{Delete}(\eta) \\ 0 & \text{otherwise.} \end{cases}$$

## Lifted Zanella process

A random walk on  $\mathcal{M}_n \times \{+1, -1\}$  with correct marginal:

For  $\gamma \in \mathcal{M}_n$ ,

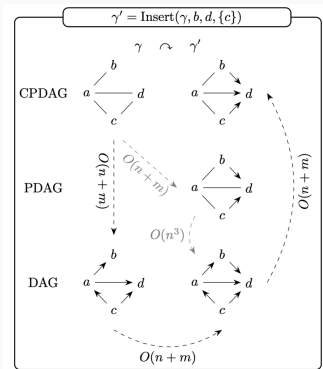
$$\lambda(\gamma^{+1} \curvearrowright \eta^{+1}) = \begin{cases} g\left(\frac{\pi\{\eta\}}{\pi\{\gamma\}}\right) & \text{if } \eta \in \text{Insert}(\gamma) \\ 0 & \text{otherwise.} \end{cases}$$

$$\lambda(\eta^{-1} \curvearrowright \gamma^{-1}) = \begin{cases} g\left(\frac{\pi\{\gamma\}}{\pi\{\eta\}}\right) & \text{if } \gamma \in \text{Delete}(\eta) \\ 0 & \text{otherwise.} \end{cases}$$

and for  $\gamma \in \mathcal{M}_n$  and  $d \in \{-1, +1\}$ ,

$$\lambda(\gamma^d \curvearrowright \gamma^{-d}) = \left( -\sum_{\eta} \lambda(\gamma^d \curvearrowright \eta^d) + \sum_{\eta} \lambda(\gamma^{-d} \curvearrowright \eta^{-d}) \right)^+.$$

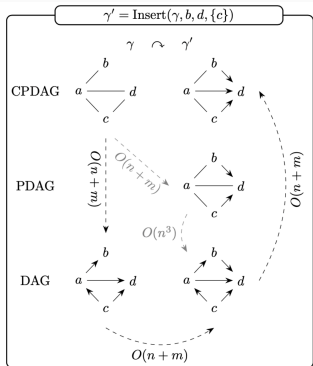
# Moving efficiently



Linear-time approach for applying a GES operator.

Previous approaches: add the inserted edge to the initial CPDAG, obtaining a PDAG associated with the new MEC  $\gamma'$ .

# Moving efficiently



Linear-time approach for applying a GES operator.

Our approach: find a consistent DAG extension of the initial CPDAG in time  $O(n + m)$ , which has the property that applying the operator directly yields a DAG from  $\gamma'$ .

## What else is there?

---

- Plug and play: `CausalInference.jl`
- Intriguing connection to the GES algorithms (greedy search for the MEC which maximises score).
- Some ideas how to handle unobserved confounders.

# What causal models does the ADNI data suggest?

