

Structure Learning for Directed Trees

Jakobsen, M. E., Shah, R. D., Bühlmann, P., Peters, J. (2022)

Jieun Shin

December 7, 2022

1. Introduction
2. Score-based Learning and Identifiability of Trees
3. Causal Additive Trees (CAT)
4. Hypothesis Testing
5. Simulation

1. Introduction

1. Introduction

- In structural causal models, one assumes that there are (causal) functions f_1, \dots, f_p such that for all

$$1 \leq i \leq p: \quad X_i := f_i(X_{\text{PA}(i)}, N_i),$$

for subsets $\text{PA}(i) \subset \{1, \dots, p\}$ and jointly independent noise variable $N = (N_1, \dots, N_p) \sim P_N$.

- For each variable X_i , one adds directed edges from its direct causes or parents $\text{PA}(i)$ into i .
- Assumptions that guarantee identifiability of the causal graph.
 - linear additive Gaussian noise models with equal noise variance (Peters and Buhlmann, 2014).
 - linear additive non-Gaussian noise model (Shimizu et al., 2006).
 - nonlinear additive noise models (Hoyer et al., 2008; Peters et al., 2014).
 - post-nonlinear additive noise models (Zhang and Hyvarinen, 2009).
 - partially-linear additive Gaussian noise models (Rothenhausler et al., 2018).
 - discrete models (Peters et al., 2011).

1. Introduction

Structure learning methods

- **Score-based learning** starts with a function ℓ assigning a population score to causal structures. Depending on the assumed model class, this function is minimized by the true structure.
- For example, when considering directed acyclic graph (DAG), the true causal DAG \mathcal{G} satisfy

$$\mathcal{G} \in \underset{\tilde{\mathcal{G}}: \tilde{\mathcal{G}} \text{ is a DAG}}{\operatorname{argmin}} \ell(\tilde{\mathcal{G}}). \quad (1)$$

- **Constraint-based learning** test for conditional independences in P_X and use these results to infer the causal structure.
- After finding the skeleton first, give it a direction.
- studied for polytrees.

1. Introduction

polytree in constraint-based learning

- Allow for multiple root as well as nodes with multiple parents.
- Chow and Liu (1968) showed how to recover an undirected Markov tree from a given discrete joint pdf using the maximum weight spanning tree (MWST) algorithm.
- Rebane and Pearl (1987) propose a method to identify the causal basins, structures constructed from nodes with at least two parents.

1. Introduction

polytree in constraint-based learning (Cont.)

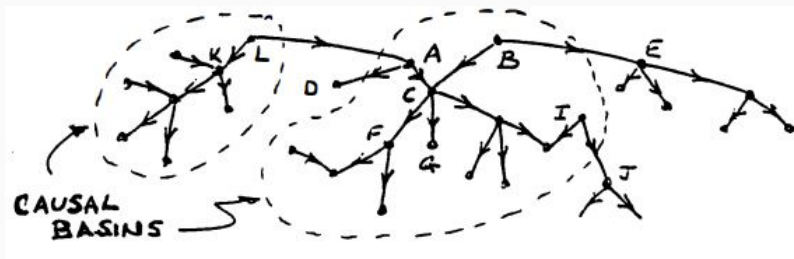


Figure 1: Example for basins Rebane and Pearl (1987).

2. Score-based Learning and Identifiability of Trees

2-1. Identifiability of Causal Additive Tree Models

graph terminology

- A directed graph $\mathcal{G} = (V, \mathcal{E})$ consists of $p \in \mathbb{N}_{>0}$ vertices (or nodes) $V = \{1, \dots, p\}$ and a collection of directed edges $\mathcal{E} \subset \{(i \rightarrow j) \equiv (i, j) : i, j \in V, i \neq j\}$
- The unique node of a directed tree \mathcal{G} with no parents is called the root node and is denoted by $\text{rt}(\mathcal{G})$.
- let \mathcal{T}_p denote the set of directed trees over $p \in \mathbb{N}_{>0}$ nodes.

2-1. Identifiability of Causal Additive Tree Models

Any tuple $(\mathcal{G}, (f_i), P_N) \in \mathcal{T}_p \times \mathcal{M}^p \times \mathcal{P}^p$ induces a structural causal model over $X = (X_1, \dots, X_p)$ given by the following structural assignments

$$X_i := f_i(X_{\text{pa}^{\mathcal{G}}(i)}) + N_i, \quad \forall 1 \leq i \leq p,$$

where $f_{\text{rt}(\mathcal{G})} \equiv 0$ and $N = (N_1, \dots, N_p) \sim P_N$, which we call a structural causal additive tree model.

- \mathcal{M} denotes all measurable functions,
- \mathcal{D}_k denotes the set of all $k \in \mathbb{N}$ times differentiable functions.
- \mathcal{P} denotes the set of mean zero probability measures on \mathbb{R} .
- For any set \mathcal{P} of probability measures, \mathcal{P}^p denotes all p -dimensional product measures on \mathbb{R}^p with marginals in \mathcal{P} .

2-1. Identifiability of Causal Additive Tree Models

- \mathcal{C}_k denotes the k times continuously differentiable functions.
- $\mathcal{P}_+ \subset \mathcal{P}$ denotes the subset for which a density is strictly positive.
- For any class $\mathcal{F} \subseteq \{f: \mathbb{R} \rightarrow \mathbb{R}\}$, $\mathcal{P}_+ \subset \mathcal{P}$ denotes the subset with a density function in \mathcal{F} .
- As a special case, we let $\mathcal{P}_G \subset \mathcal{P}_{+\mathcal{C}_\infty} := \mathcal{P}_+ \cap \mathcal{P}_{\mathcal{C}_\infty}$ denotes the subset of Gaussian probability measures.

Definition (Restricted structural causal additive tree models)

The collection of restricted structural additive tree models

$\Theta_R \subset \mathcal{T}_p \times \mathcal{D}_3^p \times \mathcal{P}_{+\mathcal{C}_3}^p$ is given by all models $\theta = (\mathcal{G}, (f_i), P_N) \in \mathcal{T}_p \times \mathcal{D}_3^p \times \mathcal{P}_{+\mathcal{C}_3}^p$ satisfying the following conditions for all $i \in \{1, \dots, p\} \setminus \{\text{rt}(\mathcal{G})\}$:

1. f_i is nowhere constant, i.e., it is not constant on any non-empty open set
2. the induced log-density ξ of $X_{\text{pa}^{\mathcal{G}}(i)}$, noise log-density ν of N_i and causal function f_i are such that there exists $x, y \in \mathbb{R}$ with $\nu''(y - f_i(x))f_i'(x) \neq 0$ such that

$$\xi''' \neq \xi'' \left(\frac{f_i''}{f_i'} - \frac{\nu''' f_i'}{\nu''} \right) - 2\nu'' f_i'' f_i' + \nu' f_i''' + \frac{\nu' \nu''' f_i'' f_i'}{\nu''} - \frac{\nu' (f_i''')^2}{f_i'},$$

2-1. Identifiability of Causal Additive Tree Models

For example, consider causal additive tree model with Gaussian noise. Then, the second condition is that f_i is not linear.

To show this, assume that

- (1) of Definition 2 is satisfied (let f_i be nowhere constant)
- (2) of Definition 2 is not satisfied (f_i is linear function.)

The log density of N_i , for all $i \in \{1, \dots, p\}$ is given by

$$\nu_i(x) = -\frac{1}{2} \log(2\pi\sigma_i^2) - \frac{x^2}{\sigma_i^2}, \quad \nu_i'(x) = \frac{x}{\sigma_i^2}, \quad \nu_i''(x) = \frac{1}{\sigma_i^2}, \quad \nu_i'''(x) = 0.$$

And there exists an $i \in \{1, \dots, p\} \setminus \{\text{rt}(\mathcal{G})\}$ such that for all

$$(x, y) := \{(x, y) \in \mathbb{R}^2 : \nu_i''(y - f_i(x)) \dot{f}_i'(x) \neq 0\} \\ \{(x, y) \in \mathbb{R}^2 : \dot{f}_i'(x) \neq 0\},$$

it holds that

$$\xi'''(x) - \xi''(x) \frac{\dot{f}_i''(x)}{\dot{f}_i'(x)} - \frac{2\dot{f}_i''(x)\dot{f}_i'(x)}{\sigma^2} = -\frac{y - f_i(x)}{\sigma^2} \left(\dot{f}_i'''(x) - \frac{(\dot{f}_i''(x))^2}{\dot{f}_i'(x)} \right). \quad (2)$$

2-1. Identifiability of Causal Additive Tree Models

Then, it hold that

$$f_i'''(x) - \frac{(f_i''(x))^2}{f_i'(x)} = \frac{\frac{\partial f''(x)}{\partial x} f'(x) - f'(x) \frac{\partial f'(x)}{\partial x}}{(f'(x))^2} = \frac{\partial}{\partial x} \left(\frac{f''(x)}{f'(x)} \right) = 0,$$

i.e., $f''(x)/f'(x)$ is constant.

Thus, the differential equation holds that

$$0 = \xi'''(x) - \xi''(x) \frac{f_i''(x)}{f_i'(x)} - \frac{2f_i''(x)f_i'(x)}{\sigma^2} = \frac{\partial}{\partial x} \left(\frac{\xi''(x)}{f'(x)} \right) - 2 \frac{f''(x)}{\sigma^2},$$

by the division with $f'(x)$. By integration this implies that

$0 = \xi''(x)/f'(x) - 2f'(x)/\sigma^2 + c_3$ such that

$\xi''(x) = 2 \exp(2c_1x + 2c_2)/\sigma^2 - c_3 \exp(c_1x + c_2)$ and

$\xi'(x) = \exp(2c_1x + 2c_2)/c_1\sigma^2 - c_e \exp(c_1x + c_2)/c_1 + c_4$ and

$$\xi(x) = \frac{\exp(2c_1x + 2c_2)}{2c_1^2\sigma^2} - \frac{c_3 \exp(c_1x + c_2)}{c_1^2} + c_4x + c_5.$$

We see that $\xi \rightarrow \infty \iff p_{X_{\text{pa}\mathcal{G}(j)}} \rightarrow \infty$ as $x \rightarrow (c_1) \cdot \infty$. it is contradiction. ■

2-1. Identifiability of Causal Additive Tree Models

Existing identifiability results for causal graphs in restricted SCM (Hoyer et al., 2008, Peters et al., 2014), for all $\theta \in \Theta_R$ and $\tilde{\theta} \in \tilde{\Theta}_R$ prove that

$$\mathcal{G} \neq \tilde{\mathcal{G}} \Rightarrow \mathcal{L}(X_\theta) \neq \mathcal{L}(X_{\tilde{\theta}}),$$

where \mathcal{L} denotes the distribution of a random variable.

The proposition is extended to result that does not assume that $\tilde{\theta}$ is a restricted causal model as below.

Proposition (Identifiability of causal additive tree models)

Suppose that $\theta \in \Theta_R$ and $\tilde{\theta} \in (\mathcal{T}, \mathcal{D}_1^p, \mathcal{P}_{\mathcal{C}_0}^p)$ respectively, it holds that

$$\mathcal{L}(X_\theta) = \mathcal{L}(X_{\tilde{\theta}}) \Rightarrow \mathcal{G} = \tilde{\mathcal{G}}.$$

2-2. Score functions

- For $\forall i \neq j$, assume that $X_i - \mathbb{E}[X_i|X_j]$ has a density.
- For any graph $\tilde{\mathcal{G}} \in \mathcal{T}_p$, we define for each node $i \in V$, Gaussian score as

$$\ell_G(\tilde{\mathcal{G}}) = \sum_{i=1}^p \ell_G(\tilde{\mathcal{G}}, i),$$

where

$$\ell_G(\tilde{\mathcal{G}}, i) := \log \left(\text{Var} \left(X_i - \mathbb{E} \left[X_i \mid X_{\text{pa}_{\tilde{\mathcal{G}}}(i)} \right] \right) \right) / 2,$$

- The Gaussian score is seen as

$$\ell_G(\tilde{\mathcal{G}}) = \inf_{Q \in \{\tilde{\mathcal{G}}\} \times \mathcal{D}_1^p \times \mathcal{P}_G^p} h(P_X, Q) - p \log(\sqrt{2\pi e}), \quad (3)$$

where $Q \in \mathcal{T}_p \times \mathcal{M}^p \times \mathcal{P}^p$ and h is cross-entropy.

2-2. Score functions

proof of (3)

Consider an SCM $\tilde{\theta} \in \{\tilde{\mathcal{G}}\} \times \mathcal{D}_1^p \times \mathcal{P}_G^p$ and $Q_{\tilde{\theta}}$ is induced distribution.

The entropy between P_X and $Q_{\tilde{\theta}}$ is then given by

$$\begin{aligned} h(P_X, Q_{\tilde{\theta}}) &:= \mathbb{E}[-\log(q_{\tilde{\theta}}(X))] \\ &= \sum_{j=1}^p \mathbb{E}[-\log(q_{\tilde{\theta}}(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}})))] \\ &= \sum_{j=1}^p h(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}), \tilde{N}_i) \end{aligned}$$

As $\tilde{N} \sim \mathcal{N}(0, \tilde{\sigma}_i^2)$ for some $\tilde{\sigma}_i > 0$,

$$\begin{aligned} h(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}), \tilde{N}_i) &= \mathbb{E} \left[-\log \left(\frac{1}{\sqrt{2\pi}\sigma_i} \exp \left(-\frac{(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}))^2}{2\tilde{\sigma}_i^2} \right) \right) \right] \\ &= \log(\sqrt{2\pi}\sigma_i) + \frac{\mathbb{E} \left((X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}))^2 \right)}{2\tilde{\sigma}_i^2}. \end{aligned}$$

2-2. Score functions

proof of (3) (Cont.)

We thus have

$$\inf_{\tilde{\sigma}_i > 0} \left\{ \log(\sqrt{2\pi}\sigma_i) + \frac{\mathbb{E} \left(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}) \right)^2}{2\tilde{\sigma}_i^2} \right\} \\ = \log(\sqrt{2\pi}) + \frac{1}{2} \log \mathbb{E} \left(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}) \right)^2 + \frac{1}{2}.$$

We conclude that

$$\inf_{Q \in \mathcal{T}_p \times \mathcal{M}^p \times \mathcal{P}^p} h(P_X, Q) = p \log(\sqrt{2\pi}) + \frac{p}{2} + \sum_{i=1}^p \frac{1}{2} \log \left(\inf_{\tilde{f}_i \in \mathcal{D}_1} \mathbb{E} \left(X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}) \right)^2 \right).$$

Finally, we have that

$$\inf_{\tilde{f}_i \in \mathcal{D}_1} \mathbb{E} \left[X_i - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}) \right]^2 = \mathbb{E} \left[\left(X_i - \mathbb{E}[X_i | X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}] \right)^2 \right] \\ + \inf_{\tilde{f}_i \in \mathcal{D}_1} \mathbb{E} \left[\left(\mathbb{E}[X_i | X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}] - \tilde{f}(X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}) \right)^2 \right] \\ = \mathbb{E} \left[\left(X_i - \mathbb{E}[X_i | X_{\text{PA}^{\tilde{\mathcal{G}}(i)}}] \right)^2 \right]. \quad \blacksquare$$

2-2. Score functions

- To choosing the best scoring graph, we use the identifiability gap between $\tilde{\mathcal{G}}$ and causal graph \mathcal{G} formed as

$$\operatorname{argmin}_{\mathcal{G} \in \mathcal{T} \setminus \{\mathcal{G}\}} \ell(\tilde{\mathcal{G}}) - \ell(\mathcal{G}). \quad (4)$$

- In Gaussian noise, the Gaussian score gap equals to

$$\ell(\tilde{\mathcal{G}}) - \ell(\mathcal{G}) = \inf_{Q \in \{\tilde{\mathcal{G}}\} \times \mathcal{D}_1^p \times \mathcal{P}_G^p} h(P_X, Q) - h(P_X) = \inf_{Q \in \{\tilde{\mathcal{G}}\} \times \mathcal{D}_1^p \times \mathcal{P}_G^p} D_{\text{KL}}(P_X || Q),$$

where D_{KL} denotes the Kullback-Leibler divergence. And it implies that

$$D_{\text{KL}}(P_X || Q) > 0.$$

- From assumption the gap is strictly positive, the score functions can be used to identify the true causal graph of a restricted structural model:

$$\mathcal{G} = \operatorname{argmin}_{\tilde{\mathcal{G}} \in \mathcal{T}_p} \ell(\tilde{\mathcal{G}}). \quad (5)$$

3. Causal Additive Trees (CAT)

Chu-Liu-Edmonds' Algorithm

- Based on a maximum spanning tree algorithm
 1. Build a complete graph with directed and weighted edges.
 2. Keep only incoming edges with the maximum scores.
 3. If there is no cycle, go to #5.
 4. If there is a cycle, pretend vertices in the cycle as one vertex and update scores for all incoming edges to the cycle; goto #2.
 5. Break all cycles by removing inappropriate edges in the cycle.

Figure 2: Chu-Liu-Edmonds' algorithm

3. Algorithm

Chu-Liu-Edmonds' Algorithm

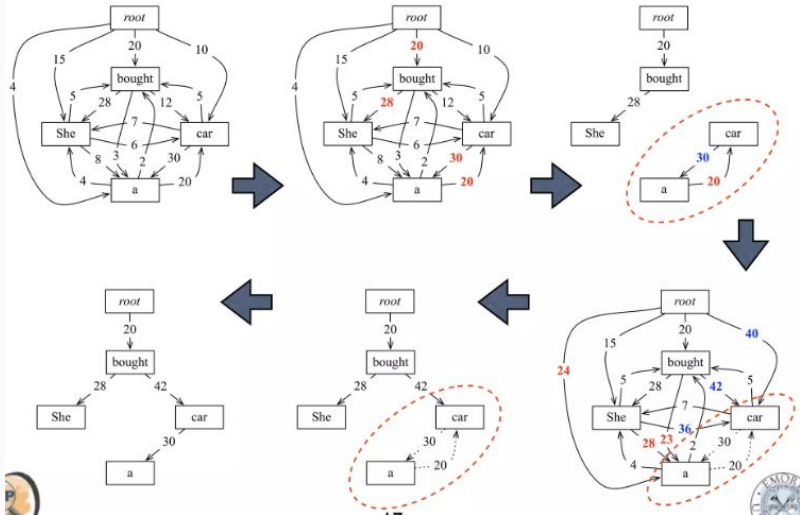


Figure 3: Chu-Liu-Edmonds' algorithm

3. Algorithm

- \mathcal{T}_p denotes directed spanning trees that is recovered by Chu-Liu-Edmonds' algorithm.
- Minimizing the Gaussian score is equivalent to minimizing a translated version of the Gaussian score function:

$$\operatorname{argmin}_{\tilde{\mathcal{G}} \in \mathcal{T}_p} \ell_G(\tilde{\mathcal{G}}) = \operatorname{argmin}_{\tilde{\mathcal{G}} \in \mathcal{T}_p} \sum_{i=1}^p \frac{1}{2} \log(\operatorname{Var}(X_i - \mathbb{E}[X_i \mid X_{\operatorname{PA} \tilde{\mathcal{G}}(i)}])) - \sum_{i=1}^p \frac{1}{2} \log(\operatorname{Var}(X_i)) \quad (6)$$

$$= \operatorname{argmin}_{\tilde{\mathcal{G}} \in \mathcal{T}_p} \sum_{i=1}^p \frac{1}{2} \log \left(\frac{\operatorname{Var}(X_i - \mathbb{E}[X_i \mid X_{\operatorname{PA} \tilde{\mathcal{G}}(i)}])}{\operatorname{Var}(X_i)} \right). \quad (7)$$

- Then, from (7), the Gaussian edge weights $w^G := (w_{ji}^G)_{j \neq i}$

$$w_{ji}^G := \frac{1}{2} \log \left(\frac{\operatorname{Var}(X_i - \mathbb{E}[X_i \mid X_j])}{\operatorname{Var}(X_i)} \right). \quad (8)$$

- Solving the problem (5) is equal to

$$\mathcal{G} = \operatorname{argmin}_{\tilde{\mathcal{G}} = (V, \tilde{\mathcal{E}}) \in \mathcal{T}_p} \sum_{(j \rightarrow i) \in \tilde{\mathcal{E}}} w_{ji}^G. \quad (9)$$

3. Algorithm

- Given the observed nodes $X = (X_1, \dots, X_p)$, we estimate the edge weights by simple plug-in estimators.
- Let us denote the conditional expectation function and its estimated by

$$\varphi_{ji}(x) := \mathbb{E}[X_i|X_j = x], \quad \hat{\varphi}_{ji}(x) := \hat{\mathbb{E}}[X_i|X_j = x] \quad (10)$$

for all $j \neq i$.

- The empirical Gaussian edge weights $\hat{w}^G = (\hat{w}_{ji}^G)_{j \neq i}$ are then given by

$$\hat{w}_{ji}^G := \frac{1}{2} \log \left(\frac{\widehat{\text{Var}}(X_i - \hat{\varphi}_{ji}(X_j))}{\widehat{\text{Var}}(X_i)} \right), \quad (11)$$

for all $j \neq i$, where $\widehat{\text{Var}}$ denotes a variance estimator.

Algorithm Causal additive trees (CAT)

- 1: **procedure** CAT (X_n , regression method)
 - 2: Run regression method to obtain $\hat{\varphi}_{ji}$ for all $j \neq i$.
 - 3: Compute empirical edge weights \hat{w}^G , see Equation (11).
 - 4: Apply Chu-Liu-Edmonds' algorithm to find MWDST with respect to \hat{w}^G .
 - 5: **return** MWDST $\hat{\mathcal{G}}$.
 - 6: **end procedure**
-

3. Algorithm

Theorem (Pointwise consistency)

Suppose that for all $j \neq i$ the following two conditions hold:

(a) if $(j \rightarrow i) \in \mathcal{E}$, $\mathbb{E} [(\hat{\varphi}_{ji}(X_j) - \varphi_{ji}(X_j))^2 \mid \tilde{\mathbf{X}}_n] \xrightarrow{P} 0$.

(b) if $(j \rightarrow i) \notin \mathcal{E}$, $\mathbb{E} [(\hat{\varphi}_{ji}(X_j) - \tilde{\varphi}_{ji}(X_j))^2 \mid \tilde{\mathbf{X}}_n] \xrightarrow{P} 0$ for some fixed $\tilde{\varphi}_{ji} : \mathbb{R} \rightarrow \mathbb{R}$.

where φ_{ji} and $\hat{\varphi}_{ji}$ are defined in Equation (7). Furthermore, suppose that Assumption 1 holds. In the large sample limit, we recover the causal graph with probability one, that is

$$P(\hat{\mathcal{G}} = \mathcal{G}) \rightarrow_n 1,$$

where $\hat{\mathcal{G}}$ is the output of Algorithm 1 using weights \hat{w}^G given by Equation (9).

4. Hypothesis Testing

4-1. Confidence Region for the Causal Tree

- The joint distribution of the estimated edge weights should be asymptotically Gaussian:
 - For all $k \in \{1, \dots, n\}$, let the vectors of squared residuals be given by

$$\hat{M}_k := \{(X_{k,i} - \hat{\varphi}(X_{k,j}))^2\}_{j \neq i} \in \mathbb{R}^{p(p-1)}$$

and let squared centered observations be given by

$$\hat{V}_k = \left\{ \left(X_{k,i} - \frac{1}{n} \sum_{m=1}^n X_{m,i} \right)^2 \right\}_{1 \leq i \leq p} \in \mathbb{R}^p.$$

- Further let

$$\hat{\mu} := \frac{1}{n} \sum_{k=1}^n \hat{M}_k, \quad \hat{\nu} := \frac{1}{n} \sum_{k=1}^n \hat{V}_k.$$

- The empirical Gaussian edge weight for $j \rightarrow i$ is given by $\log(\hat{\mu}_{ji}/\hat{\nu}_i)/2$.
- Let $\hat{\Sigma}_M \in \mathbb{R}^{p(p-1) \times p(p-1)}$ denote the empirical variances of \hat{M}_k .
- Let $\hat{\Sigma}_V \in \mathbb{R}^{p \times p}$ denote the empirical variances of \hat{V}_k .
- Let $\hat{\Sigma}_{MV} \in \mathbb{R}^{p(p-1) \times p}$ denote the empirical covariance of \hat{M}_k and \hat{V}_k .
-

$$\hat{\Sigma} := \begin{pmatrix} \hat{\Sigma}_M & \hat{\Sigma}_{MV} \\ \hat{\Sigma}_{MV} & \hat{\Sigma}_V \end{pmatrix} := \begin{pmatrix} \hat{M}_k \hat{M}_k^T - \hat{\mu} \hat{\mu}^T & \hat{M}_k \hat{V}_k^T - \hat{\mu} \hat{\nu}^T \\ \hat{V}_k \hat{M}_k^T - \hat{\nu} \hat{\mu}^T & \hat{V}_k \hat{V}_k^T - \hat{\nu} \hat{\nu}^T \end{pmatrix}$$

4-1. Confidence Region for the Causal Tree

- Set the confidence interval

$$\hat{u}_{ji}, \hat{l}_{ji} := \frac{1}{2} \log \left(\frac{\hat{\mu}_{ji}}{\hat{\nu}_i} \right) \pm z_\alpha \frac{\hat{\sigma}_{ji}}{2\sqrt{n}} = \hat{w}_{ji}^G \pm z_\alpha \frac{\hat{\sigma}_{ji}}{2\sqrt{n}}, \quad (12)$$

where $\hat{\sigma}_{ji}^2 := \frac{\hat{\Sigma}_{M,ji,ji}}{\hat{\mu}_{ji}^2} + \frac{\hat{\Sigma}_{V,i,i}}{\hat{\nu}_i^2} - \frac{\hat{\Sigma}_{MV,ji,i}}{\hat{\mu}_{ji}\hat{\nu}_i}$.

By Bonferroni correction, z_α denotes the upper $\alpha/\{2p(p-1)\}$ quantile of a standard normal distribution.

- The region of directed trees formed of minimizers of the score with edge weights in the confidence hyper-rectangle:

$$\hat{C}_{\text{Bon}} := \hat{C}(\hat{l}, \hat{u}) := \left\{ \underset{\tilde{G}=(V,\tilde{E}) \in \tilde{T}_p}{\operatorname{argmin}} \sum_{(j \rightarrow i)} w'_{ji} : \forall j \neq i, w'_{ji} \in [\hat{l}_{ji}, \hat{u}_{ji}] \right\}.$$

4-1. Confidence Region for the Causal Tree

Theorem (Confidence region)

Suppose the following conditions hold:

1. there exists $\xi > 0$ such that $\mathbb{E}|X|^{4+\xi} < \infty$
2. there exists $\xi > 0$ such that for all $j \neq i$, $\mathbb{E}[|\hat{\varphi}_{ji}(X_j) - \varphi_{ji}(X_j)|^{4+\xi} | \tilde{\mathbf{X}}_n] = O_p(1)$
3. $\text{Var} \left((\hat{M}_1^T, \hat{V}_1^T)^T | \tilde{\mathbf{X}}_n \right) \xrightarrow{P} {}_n\Sigma$, where Σ is constant with strictly positive diagonal
4. for $(j \rightarrow i) \in \mathcal{E}$, $\sqrt{n} \mathbb{E} [(\hat{\varphi}_{ji}(X_{k,j}) - \varphi_{ji}(X_{k,j}))^2 | \tilde{\mathbf{X}}_n] \xrightarrow{P} {}_n0$

Then

$$\liminf_{n \rightarrow \infty} P(\mathcal{G} \in \hat{\mathcal{C}}_{\text{Bon}}) \geq 1 - \alpha$$

4-2. Testing of Substructures

But, the confidence region \hat{C}_{Bon} will not be possible to compute it in practice..
(due to the ranges of w'_{ji})

A substructure restriction $\mathcal{R} = (\mathcal{E}_{\mathcal{R}}, \mathcal{E}_{\mathcal{R}}^{\text{miss}}, r)$ exists on the node V , where

- $\mathcal{E}_{\mathcal{R}}$ denotes existing edges
- $\mathcal{E}_{\mathcal{R}}^{\text{miss}}$ denotes missing edges
- r denotes a specific root node r

The null hypothesis for testing substructure:

$$\mathcal{H}_0(\mathcal{R}) : \mathcal{E}_{\mathcal{R}} \setminus \mathcal{E} = \emptyset, \mathcal{E} \setminus \mathcal{E}_{\mathcal{R}}^{\text{miss}} = \mathcal{E}, r = \text{rt}(\mathcal{G}), \quad (13)$$

where $\mathcal{G} = (V, \mathcal{E})$ is the true graph.

4-2. Testing of Substructures

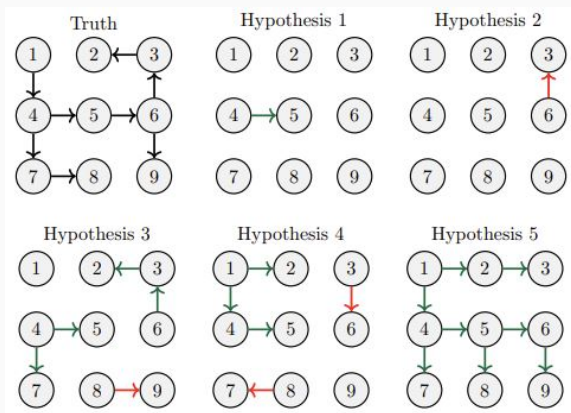


Figure 4: Illustration for hypothesis testing (green is $\mathcal{E}_{\mathcal{R}}$, red is $\mathcal{E}_{\mathcal{R}}^{\text{miss}}$).

- Hypothesis 1 consists of the restriction $\mathcal{R} = \mathcal{E}_{\mathcal{R}}$, where $\mathcal{E}_{\mathcal{R}} := \{(X_4 \rightarrow X_5)\}$. (*true*)
- Hypothesis 2 consists of the restriction $\mathcal{R} = \mathcal{E}_{\mathcal{R}}^{\text{miss}}$, where $\mathcal{E}_{\mathcal{R}}^{\text{miss}} := \{(X_6 \rightarrow X_3)\}$. (*false*)

4-2. Testing of Substructures

1. CheckC (Check confidence interval)

- Let $\mathcal{T}_p(\mathcal{R}) \subset \mathcal{T}_p$ be the set of all directed trees satisfying the substructure restriction \mathcal{R} .
- Suppose that the causal directed tree \mathcal{G} satisfies \mathcal{R} , i.e., $\mathcal{G} \in \mathcal{T}_p(\mathcal{R})$.
- From Chu-Liu-Edmond's algorithm, there exists a graph in $w' = (w'_{ji})_{j \neq i}$, with $\hat{l}_{ji} \leq w' \leq \hat{u}_{ji}$ for all $j \neq i$.
- Find $S_{\mathcal{T}_p(\mathcal{R})}(\hat{l})$ and $\mathcal{G}_{\mathcal{T}_p(\mathcal{R})}^*(\hat{l})$. Hence, it must hold that

$$S_{\mathcal{T}_p(\mathcal{R})}(\hat{l}) \leq S_{\mathcal{T}_p(\mathcal{R})}(w') = S_{\mathcal{T}_p}(w') \leq S_{\mathcal{T}_p}(\hat{u}).$$

- Define CheckC test function as

$$\psi_{\mathcal{R}}^{\text{CheckC}} := \begin{cases} 0, & \text{if } S_{\mathcal{T}_p(\mathcal{R})}(\hat{l}) \leq S_{\mathcal{T}_p(\mathcal{R})}(\hat{u}), \\ 1, & \text{otherwise.} \end{cases} \quad (14)$$

4-2. Testing of Substructures

Algorithm Hypothesis testing of $H_0(\mathcal{R})$ using the CheckC test

- 1: **procedure** CheckC ($\mathcal{R} = (\mathcal{E}_{\mathcal{R}}, \mathcal{E}_{\mathcal{R}}^{\text{miss}}, r)$, $\hat{l} = (\hat{l}_{ji})_{j \neq i}$, $\hat{u} = (\hat{u}_{ji})_{j \neq i}$).
 - 2: Initialize fully connected graph $\mathcal{H} := \{(j \rightarrow i) : i, j \in V, j \neq i\}$.
 - 3: For each $(j \rightarrow i) \in \mathcal{E}_{\mathcal{R}}$, delete from \mathcal{H} the edges
 $\{(k \rightarrow i) : k \in V \setminus \{j\}\} \cup \{i \rightarrow j\}$.
 - 4: For each $(j \rightarrow i) \in \mathcal{E}_{\mathcal{R}}^{\text{miss}}$ *from the edges* $\{(j \rightarrow i)$.
 - 5: If root $r \in \mathcal{R}$, delete from \mathcal{H} the edges $\{(j \rightarrow r) : j \in V\}$.
 - 6: Apply Chu-Liu-Edmonds' algorithm to find $S_{\mathcal{T}_p(\mathcal{R})}(\hat{l})$ and $\mathcal{G}_{\mathcal{T}_p(\mathcal{R})}^*(\hat{l})$,
 the minimum \hat{u} -weighted directed spanning subtree of \mathcal{H} .
 - 7: If $S_{\mathcal{T}_p(\mathcal{R})}(\hat{l}) \leq S_{\mathcal{T}_p}(\hat{u})$, then set $\psi_{\mathcal{R}}^{\text{CheckC}} := 0$, otherwise set $\psi_{\mathcal{R}}^{\text{CheckC}} := 1$.
 - 8: **return** $\psi_{\mathcal{R}}^{\text{CheckC}}$.
 - 9: **end procedure**
-

4-2. Testing of Substructures

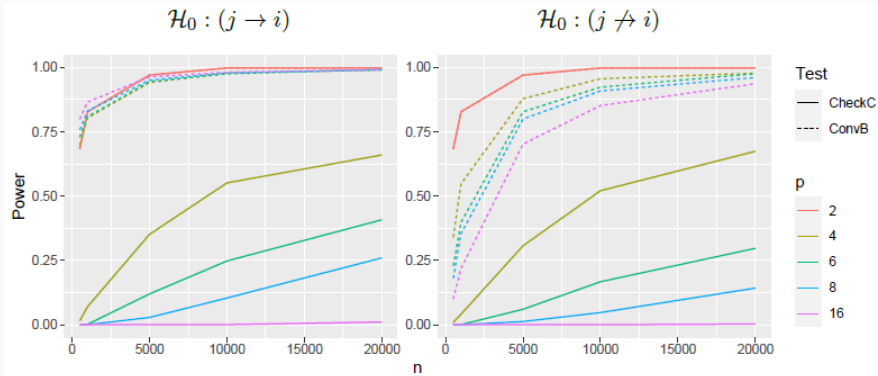


Figure 5: The power of the proposed testing procedure for simple hypotheses.

In simulation results, the power using the CheckC is lower than the alternative thing.

4-3. Converging Bounds

- Consider a true null hypothesis $H_0(\mathcal{R})$, i.e., a substructure restriction \mathcal{R} which is satisfied by the causal graph \mathcal{G} .
- Suppose that $\mathcal{G} \in \hat{\mathcal{C}}_{\text{Bon}}$, which implies exists a graph in $w' = (w'_{ji})_{j \neq i}$, with $\hat{l}_{ji} \leq w' \leq \hat{u}_{ji}$ for all $j \neq i$.
- Define the edge weights as

$$\check{w}_{ji} = \begin{cases} \hat{u}_{ji} & \text{if } [\exists k \neq j : (k \rightarrow i) \in \mathcal{E}_{\mathcal{R}}] \vee [(i \rightarrow j) \in \mathcal{E}_{\mathcal{R}}] \vee [(j \rightarrow i) \in \mathcal{E}_{\mathcal{R}}^{\text{miss}}] \vee [i = j] \\ \hat{l}_{ji} & \text{otherwise.} \end{cases}$$

- Define ConvB test function as

$$\psi_{\mathcal{R}}^{\text{ConvB}} = \begin{cases} 0, & \text{if } \mathcal{G}_{\mathcal{T}_p}^*(\check{w}) \text{ satisfies } \mathcal{R} \\ 1, & \text{otherwise.} \end{cases} \quad (15)$$

4-3. Converging Bounds

Algorithm Hypothesis testing of $\mathcal{H}_0(\mathcal{R})$ using the ConvB test

- 1: **procedure** CoNVB $\left(\mathcal{R} = (\mathcal{E}_{\mathcal{R}}, \mathcal{E}_{\mathcal{R}}^{\text{miss}}, r), \hat{l} = (\hat{l}_{ji})_{j \neq i}, \hat{u} = (\hat{u}_{ji})_{j \neq i} \right)$
 - 2: Initialize $\check{w} := \hat{l}$.
 - 3: For each $(j \rightarrow i) \in \mathcal{E}_{\mathcal{R}}$ and all $k \in V \setminus \{j\}$, set $\check{w}_{ki} := \hat{u}_{ki}$.
 - 4: For each $(j \rightarrow i) \in \mathcal{E}_{\text{Riss}}$, set $\check{w}_{ji} := \hat{u}_{ji}$.
 - 5: If root $r \in \mathcal{R}$, then for all $j \in V$, set $\check{w}_{jr} := \hat{u}_{lr}$.
 - 6: Apply Chu-Liu-Edmonds' algorithm to find $\mathcal{G}_{\mathcal{T}_p}^*(\check{w})$.
 - 7: If $\mathcal{G}_{\mathcal{T}_p}^*(\check{w})$ satisfies \mathcal{R} , then set $\psi_{\mathcal{R}}^{\text{ConvB}} := 0$, otherwise set $\psi_{\mathcal{R}}^{\text{ConvB}} := 1$.
 - 8: **return** $\psi_{\mathcal{R}}^{\text{ConvB}}$.
 - 9: **end procedure**
-

5. Simulation

5. Simulation

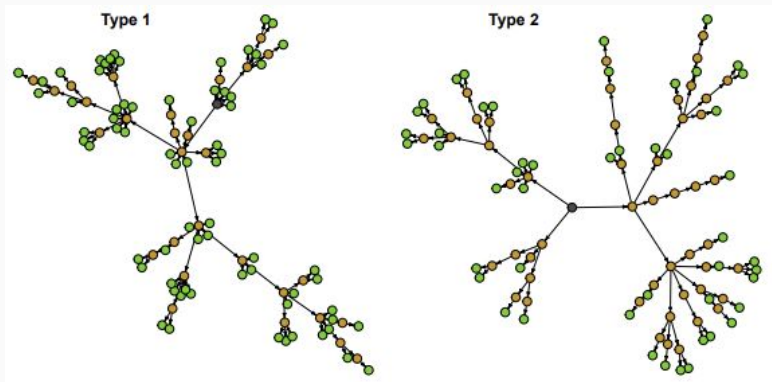


Figure 6: Illustration of Type 1 (many leaf nodes) and Type 2 (many branch nodes) directed trees over $p = 100$ nodes.

- green (leaf node), brown (branch node), black (root node)
- The Type 1 tree contains 70 leaf nodes, while the Type 2 tree only contains 49 leaf nodes.

5. Simulation

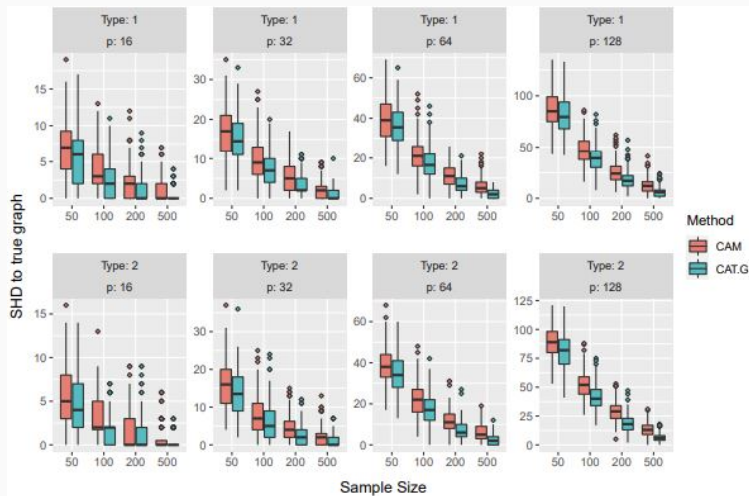


Figure 7: Causal additive tree models with Gaussian noise

Boxplots of the SHD (Structural Hamming Distance) performance of CAM (Buhlmann et al., 2014) and CAT.G