

Causal discovery: noise-based methods

Charles K. Assaad, Emilie Devijver, Eric Gaussier

charles.assaad@ens-lyon.fr

Table of content

Preliminaries

Bivariate causal discovery

Multivariate causal discovery

Conclusion

Table of content

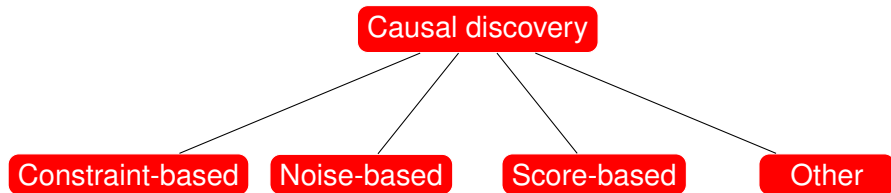
Preliminaries

Bivariate causal discovery

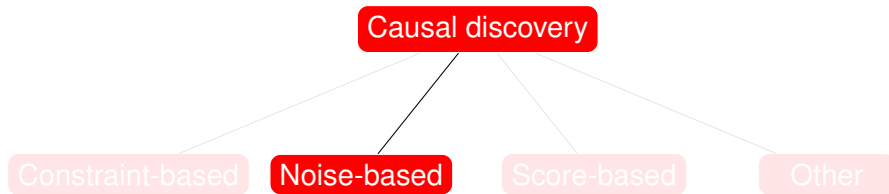
Multivariate causal discovery

Conclusion

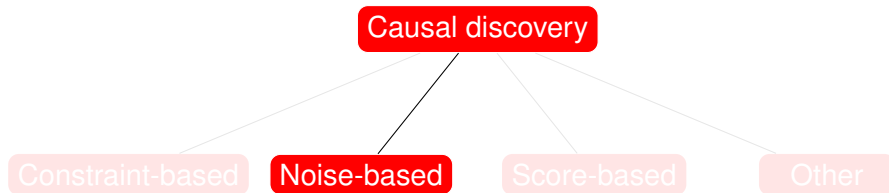
Causal discovery



Causal discovery



Causal discovery



Noise-based: find footprints in the noise that imply causal asymmetry.

Recap about causal graphical models

Causal sufficiency

$$\forall X \leftarrow Z \rightarrow Y, \text{ if } X, Y \in \mathcal{V} \text{ then } Z \in \mathcal{V}.$$

Topological ordering: Consider a causal DAG $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and a topological ordering $\mathcal{T} = \{X_1, \dots, X_p\}$. If $X_i \rightarrow X_j$ in \mathcal{G} then $i < j$.

Recap about structural causal models (1/2)

$V = \{X_1, X_2, \dots, X_n\}$ set of endogenous variables

$U = \{\zeta_1, \zeta_2, \dots, \zeta_n\}$ corresponding set of exogenous variables.

Suppose that each endogenous variable X_i is a function of its parents in V together with ζ_i :

$$X_i = f_i(\text{Parents}(X_i), \zeta_i).$$

Graphical representation is including only the endogenous variables V , and we use $\text{Parents}(X_i)$ to denote the set of endogenous parents of X_i .

Recap about structural causal models (1/2)

$V = \{X_1, X_2, \dots, X_n\}$ set of endogenous variables

$U = \{\zeta_1, \zeta_2, \dots, \zeta_n\}$ corresponding set of exogenous variables.

Suppose that each endogenous variable X_i is a function of its parents in V together with ζ_i :

$$X_i = f_i(\text{Parents}(X_i), \zeta_i).$$

Graphical representation is including only the endogenous variables V , and we use $\text{Parents}(X_i)$ to denote the set of endogenous parents of X_i .

Recap about structural causal models (1/2)

$V = \{X_1, X_2, \dots, X_n\}$ set of endogenous variables

$U = \{\zeta_1, \zeta_2, \dots, \zeta_n\}$ corresponding set of exogenous variables.

Suppose that each endogenous variable X_i is a function of its parents in V together with ζ_i :

$$X_i = f_i(\text{Parents}(X_i), \zeta_i).$$

Graphical representation is including only the endogenous variables V , and we use $\text{Parents}(X_i)$ to denote the set of endogenous parents of X_i .

Recap about structural causal models (2/2)

Independent Mechanism Principle

In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other conditional distributions.

- ▶ Independence of noises, conditional independence of structures
- ▶ Independence of information contained in mechanisms
- ▶ Intervenability, autonomy, modularity, invariance, transfer

If the system of equations is acyclic, an assignment of values to the exogenous variables $\zeta_1, \zeta_2, \dots, \zeta_n$ uniquely determines the values of all the variables in the model. Then, if we have a probability distribution P' over the values of variables in ζ , this will induce a unique probability distribution P on V .

Recap about structural causal models (2/2)

Independent Mechanism Principle

In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other conditional distributions.

- ▶ Independence of noises, conditional independence of structures
- ▶ Independence of information contained in mechanisms
- ▶ Intervenability, autonomy, modularity, invariance, transfer

If the system of equations is acyclic, an assignment of values to the exogenous variables $\zeta_1, \zeta_2, \dots, \zeta_n$ uniquely determines the values of all the variables in the model. Then, if we have a probability distribution P' over the values of variables in ζ , this will induce a unique probability distribution P on V .

The intuition behind the noise (1/2)

$$\text{Suppose } \begin{cases} X := \zeta_x \\ Y := 2X + \zeta_y \end{cases}$$

The intuition behind the noise (1/2)

$$\text{Suppose } \begin{cases} X := \zeta_x \\ Y := 2X + \zeta_y \end{cases}$$

Given $P(X, Y)$, one can detect $X - Y$ but what about orientation?

The intuition behind the noise (1/2)

$$\text{Suppose } \begin{cases} X := \zeta_x \\ Y := 2X + \zeta_y \end{cases}$$

Given $P(X, Y)$, one can detect $X - Y$ but what about orientation?

$$Y := 2X + \zeta_y ?$$

or

$$X := \frac{Y}{2} + \zeta_x ?$$

The intuition behind the noise (1/2)

$$\text{Suppose } \begin{cases} X := \xi_x \\ Y := 2X + \xi_y \end{cases}$$

Given $P(X, Y)$, one can detect $X - Y$ but what about orientation?

$$Y := 2X + \xi_y ?$$

or

$$X := \frac{Y}{2} + \xi_x ?$$

Without further assumption we cannot know.

The intuition behind the noise (1/2)

$$\text{Suppose } \begin{cases} X := \xi_x \\ Y := 2X + \xi_y \end{cases}$$

Given $P(X, Y)$, one can detect $X - Y$ but what about orientation?

$$Y := 2X + \xi_y ?$$

or

Without further assumption we cannot know.

$$X := \frac{Y}{2} + \xi_x ?$$

Assume that the noise follow a uniform distribution on $\{-1, 0, 1\}$

The intuition behind the noise (1/2)

$$\text{Suppose } \begin{cases} X := \xi_x \\ Y := 2X + \xi_y \end{cases}$$

Given $P(X, Y)$, one can detect $X - Y$ but what about orientation?

$$Y := 2X + \xi_y ?$$

or

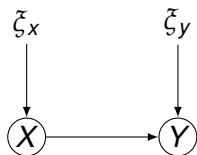
Without further assumption we cannot know.

$$X := \frac{Y}{2} + \xi_x ?$$

Assume that the noise follow a uniform distribution on $\{-1, 0, 1\}$

X	Y	$\xi_y = Y - 2X$	$\xi_x = X - Y/2$
1	2	$0 \in \{-1, 0, 1\}$	$0 \in \{-1, 0, 1\}$
3	6	$0 \in \{-1, 0, 1\}$	$0 \in \{-1, 0, 1\}$
4	9	$1 \in \{-1, 0, 1\}$	$-0.5 \notin \{-1, 0, 1\}$

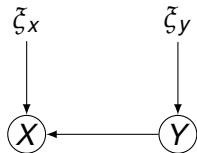
The intuition behind the noise (2/2)



$$M_1 : \begin{cases} X := f_x(\xi_x) \\ Y := f_y(X, \xi_y) \end{cases}$$

- ▶ $X \perp\!\!\!\perp_G \xi_y$
- ▶ $Y \not\perp\!\!\!\perp_G \xi_x$

Backwards model:



$$M_2 : \begin{cases} Y := g_y(\xi_y) \\ X := g_x(Y, \xi_x) \end{cases}$$

- ▶ $X \not\perp\!\!\!\perp_G \xi_y$
- ▶ $Y \perp\!\!\!\perp_G \xi_x$

Noise based question

Main question: Given $P(\mathcal{V})$ a compatible probability distribution of \mathcal{G} , can we discover \mathcal{G} ?

Noise based question

Main question: Given $P(\mathcal{V})$ a compatible probability distribution of \mathcal{G} , can we discover \mathcal{G} ? **No!**

Noise based question

Main question: Given $P(\mathcal{V})$ a compatible probability distribution of \mathcal{G} , can we discover \mathcal{G} ? **No!**

It is possible that $Y \perp\!\!\!\perp_P \tilde{\zeta}_X$.

Noise based question

Main question: Given $P(\mathcal{V})$ a compatible probability distribution of \mathcal{G} , can we discover \mathcal{G} ? **No!**

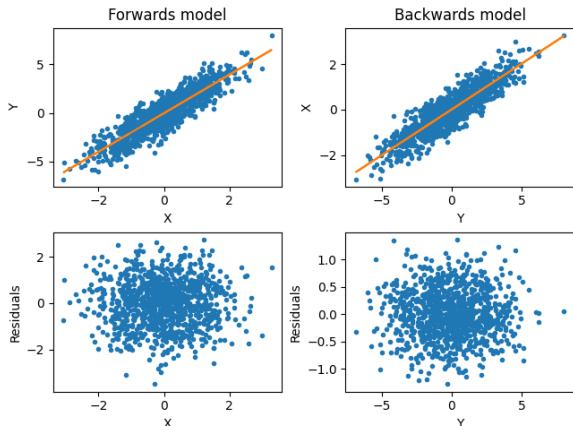
It is possible that $Y \perp\!\!\!\perp_P \zeta_X$.

Example:

$$X \sim N(0, 1)$$

$$\zeta_Y \sim N(0, 1)$$

$$Y := 2X + \zeta_Y$$



Noise based question

Main question: Given $P(\mathcal{V})$ a compatible probability distribution of \mathcal{G} , can we discover \mathcal{G} ? **No!**

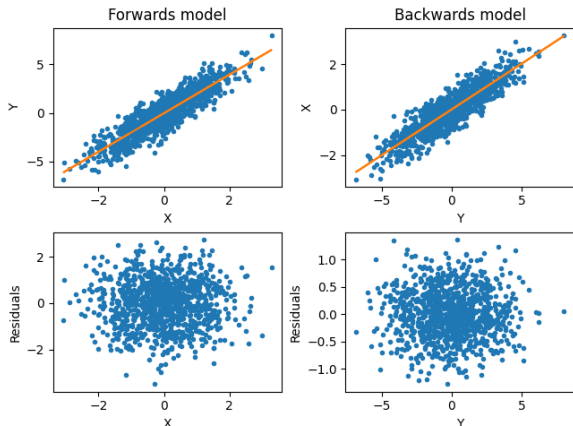
It is possible that $Y \perp\!\!\!\perp_P \zeta_X$.

Example:

$$X \sim N(0, 1)$$

$$\zeta_Y \sim N(0, 1)$$

$$Y := 2X + \zeta_Y$$



\implies The Markov equivalence class is the best we can do!

Table of content

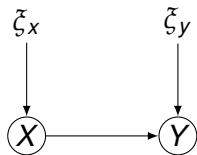
Preliminaries

Bivariate causal discovery

Multivariate causal discovery

Conclusion

The linear case (1/2)



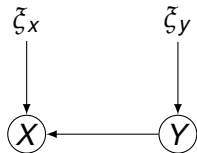
$$M_1 : \begin{cases} X := \zeta_x \\ Y := aX + \zeta_y \end{cases}$$

$$\triangleright X \perp\!\!\!\perp_G \zeta_y$$

$$\triangleright Y \not\perp\!\!\!\perp_G \zeta_x$$

When $Y \perp\!\!\!\perp_P \zeta_x$?

Backwards model:



$$M_2 : \begin{cases} Y := \zeta_y \\ X := bY + \zeta_x \end{cases}$$

$$\begin{aligned} \zeta_x &= X - bY \\ &= X - b(aX + \zeta_y) \\ &= (1 - ba)X - b\zeta_y \end{aligned}$$

The linear case (2/2)

$$Y = aX + \zeta_y$$

$$\tilde{\zeta}_x = (1 - ba)X - b\zeta_y$$

When $Y \perp\!\!\!\perp_P \tilde{\zeta}_x$?

The linear case (2/2)

$$Y = aX + \xi_y$$

$$\xi_x = (1 - ba)X - b\xi_y$$

When $Y \perp\!\!\!\perp_P \xi_x$?

Theorem (Darmois-Skitovich): Let X_1, \dots, X_n be independent, non degenerate random variables. If for two linear combinations:

$$I_1 = a_1 X_1 + \dots + a_n X_n$$

$$I_2 = b_1 X_1 + \dots + b_n X_n$$

are independent, then each X_i is normally distributed.

The linear non gaussian case (1/2)

Theorem (identifiability of linear non-Gaussian models): Assume that $P(X, Y)$ admits the linear model

$$Y := aX + \tilde{\zeta}_y, \quad X \perp\!\!\!\perp_P \tilde{\zeta}_y,$$

with continuous random variables X , $\tilde{\zeta}_y$, and Y . Then there exists $b \in \mathbb{R}$ and a random variable $\tilde{\zeta}_x$ such that

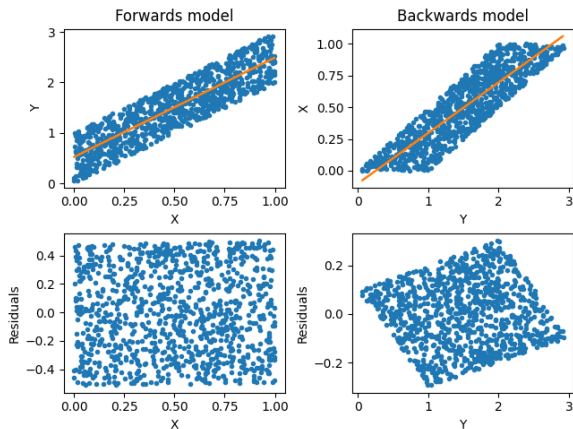
$$X := bY + \tilde{\zeta}_x, \quad Y \perp\!\!\!\perp_P \tilde{\zeta}_x,$$

if and only if $\tilde{\zeta}_y$ and X are Gaussian.
(proof on board)

The linear non gaussian case (2/2)

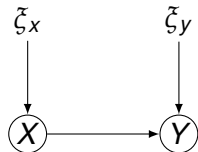
Example:

$$X \sim U(0, 1)$$
$$\xi_y \sim U(0, 1)$$
$$Y := 2X + \xi_y$$



The non linear case (1/3)

Continuous additive noise models



$$M_1 : \begin{cases} X := \tilde{\zeta}_x \\ Y := f_y(X) + \tilde{\zeta}_y \end{cases}$$

- ▶ $X \perp\!\!\!\perp_G \tilde{\zeta}_y$
- ▶ $Y \not\perp\!\!\!\perp_G \tilde{\zeta}_x$

When $Y \perp\!\!\!\perp_P \tilde{\zeta}_x$?

The non linear case (2/3)

Theorem (identifiability of additive noise models): Assume that $P(X, Y)$ admits the non-linear additive noise model

$$Y := f_Y(X) + \zeta_Y, \quad X \perp\!\!\!\perp_P \zeta_Y,$$

with continuous random variables X , ζ_Y , and Y . Then there exists $g(\cdot)$ and random variable ζ_X such that

$$X := f_X(Y) + \zeta_X, \quad Y \perp\!\!\!\perp_P \zeta_X,$$

if and only if *Complicated Condition* is satisfied.
(Hoyer et al, 2008)

The non linear case (2/3)

Theorem (identifiability of additive noise models): Assume that $P(X, Y)$ admits the non-linear additive noise model

$$Y := f_y(X) + \zeta_y, \quad X \perp\!\!\!\perp_P \zeta_y,$$

with continuous random variables X , ζ_y , and Y . Then there exists $g(\cdot)$ and random variable ζ_x such that

$$X := f_x(Y) + \zeta_x, \quad Y \perp\!\!\!\perp_P \zeta_x,$$

if and only if *Complicated Condition* is satisfied.

(Hoyer et al, 2008)

Complicated Condition: The triple $(f_y, P(X), P(\zeta_y))$ solves the following differential equation for all x, y with $(\log P(\zeta_y))''(y - f_y(x))f'(x) \neq 0$.

The non linear case (3/3)

- ▶ The space that satisfy the condition is a 3-dimensional space;
The space of continuous distributions is infinite dimensional;
⇒ we have identifiability for most distributions.
- ▶ If the noise is Gaussian, then the only functional form that satisfies Complicated Condition is linearity.
- ▶ If the function is linear and the noise is non-Gaussian, then one can't fit a linear backwards model **but** one can fit a non-linear backwards models.

Causal order discovery procedure in the bivariate case

Given $P(X, Y)$ and a dependence estimator $\hat{\lambda}$

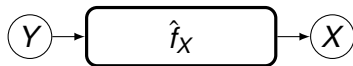
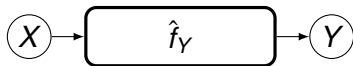
Procedure:

Causal order discovery procedure in the bivariate case

Given $P(X, Y)$ and a dependence estimator \hat{l}

Procedure:

1. Fit \hat{f}_Y and \hat{f}_X :

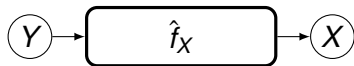
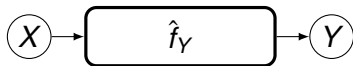


Causal order discovery procedure in the bivariate case

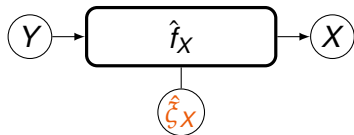
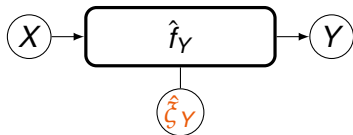
Given $P(X, Y)$ and a dependence estimator $\hat{\gamma}$

Procedure:

1. Fit \hat{f}_Y and \hat{f}_X :



2. Compute residuals $\hat{\zeta}_Y$ and $\hat{\zeta}_X$:

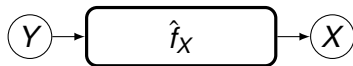
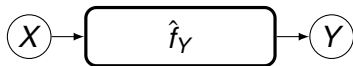


Causal order discovery procedure in the bivariate case

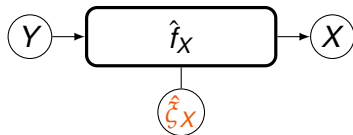
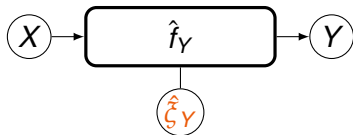
Given $P(X, Y)$ and a dependence estimator \hat{l}

Procedure:

1. Fit \hat{f}_Y and \hat{f}_X :



2. Compute residuals $\hat{\xi}_Y$ and $\hat{\xi}_X$:



3. Order:

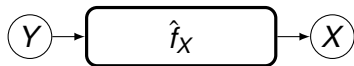
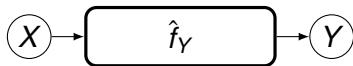
- ▶ $\mathcal{T} = [X, Y]$ if $\hat{l}(x, \hat{\xi}_Y) < \hat{l}(y, \hat{\xi}_X)$
- ▶ $\mathcal{T} = [Y, X]$ if $\hat{l}(y, \hat{\xi}_X) < \hat{l}(x, \hat{\xi}_Y)$

Causal order discovery procedure in the bivariate case

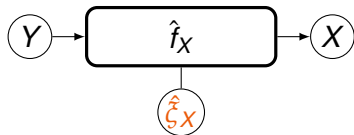
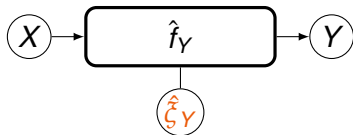
Given $P(X, Y)$ and a dependence estimator $\hat{\lambda}$

Procedure:

1. Fit \hat{f}_Y and \hat{f}_X :



2. Compute residuals $\hat{\xi}_Y$ and $\hat{\xi}_X$:



3. Order:

- ▶ $\mathcal{T} = [X, Y]$ if $\hat{\lambda}(x, \hat{\xi}_Y) < \hat{\lambda}(y, \hat{\xi}_X)$
- ▶ $\mathcal{T} = [Y, X]$ if $\hat{\lambda}(y, \hat{\xi}_X) < \hat{\lambda}(x, \hat{\xi}_Y)$

4. Output (suppose $\mathcal{T} = [X, Y]$):

- ▶ $X \rightarrow Y$ if $X \perp\!\!\!\perp_P \hat{\xi}_Y$ and $Y \not\perp\!\!\!\perp_P \hat{\xi}_X$

Table of content

Preliminaries

Bivariate causal discovery

Multivariate causal discovery

Conclusion

Minimality

Minimality condition A DAG \mathcal{G} compatible with a probability distribution P is said to satisfy the minimality condition if P is not compatible with any proper subgraph of \mathcal{G} .

Minimality

Minimality condition A DAG \mathcal{G} compatible with a probability distribution P is said to satisfy the minimality condition if P is not compatible with any proper subgraph of \mathcal{G} .

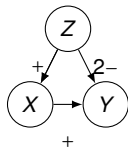
Remark: faithfulness \implies minimality.

Minimality and d-sep

Theorem (implication of minimality on d-sep): Consider the random vector \mathcal{V} and assume that the joint distribution has a density with respect to a product measure. Suppose that $P(\mathcal{V})$ is Markov with respect to \mathcal{G} . Then $P(\mathcal{V})$ satisfies the minimality condition iff $\forall X \in \mathcal{V}$ and $\forall Y \in \text{Parents}(X, \mathcal{G})$,
 $X \not\perp_P Y \mid \text{Parents}(X, \mathcal{G}) \setminus \{Y\}$.
(proof on board)

Violation of minimality

Example 1: canceling out



Example 2: constant functions

Linear non gaussian

Theorem (LiNGAM) Assume a linear SCM with graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and a compatible distribution $P(\mathcal{V})$ such that $\forall Y \in \mathcal{V}$

$$Y := \sum_{X \in \text{Parents}(Y, \mathcal{G})} a_{XY} X + \xi_Y$$

where all ξ_Y are jointly independent and non-Gaussian distributed. Additionally, we require that $\forall Y \in \mathcal{V}, X \in \text{Parents}(Y, \mathcal{G}), a_{XY} \neq 0$. Then, the graph \mathcal{G} is identifiable from $P(\mathcal{V})$.

(proof in (Shimizu et al, 2011))

The LiNGAM algorithm

Algorithm 1 LiNGAM

Input: $P(\mathcal{V})$

Output: \mathcal{G}

```
1: Form an empty graph  $\mathcal{G}$  on vertex set  $\mathcal{V} = \{X_1, \dots, X_p\}$ 
2: Let  $S = \{1, \dots, p\}$  and  $\mathcal{T} = []$ 
3: repeat
4:    $H = []$ 
5:   for  $i \in S$  do
6:     for  $j \in S \setminus \{i\}$  do
7:        $\hat{\zeta}_{ij} = X_j - \frac{\text{cov}(X_i, X_j)}{\text{var}(X_i)} X_i$ 
8:     end for
9:      $h = \sum_{j \in S \setminus \{i\}} \lambda(X_i, \hat{\zeta}_{ij})$ 
10:     $H = [H, h]$ 
11:   end for
12:    $i^* = \arg \min_{i \in S} H$ 
13:    $S = S \setminus \{i^*\}$ 
14:    $\mathcal{T} = [\mathcal{T}, i^*]$ 
15:    $\forall j \in S, X_j = \hat{\zeta}_{i^*j}$ 
16: until  $|S| = 0$ 
17: Append( $\mathcal{T}, S_0$ )
18: Construct a strictly lower triangular matrix by following the order in  $\mathcal{T}$ , and estimate the connection strengths  $a_{i,j}$  by using some conventional covariance-based regression.
19: if  $a_{i,j} > 0$  then
20:   Add  $X_i \rightarrow X_j$  to  $\mathcal{G}$ 
21: end if
22: Return  $\mathcal{G}$ 
```

Additive noise models

Theorem (ANM) Assume a non-linear SCM with graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and a compatible distribution $P(\mathcal{V})$ that satisfy the minimality condition with respect to \mathcal{G} . $\forall Y \in \mathcal{V}$

$$Y := f(\text{Parents}(Y, \mathcal{G})) + \xi_Y$$

where all ξ_Y are jointly independent. Then, the graph \mathcal{G} is identifiable from $P(\mathcal{V})$.

(proof in (Peters et al, 2014))

The ANM algorithm

Algorithm 2 ANM

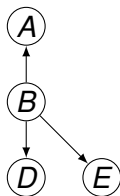
Input: $P(\mathcal{V})$

Output: \mathcal{G}

```
1: Form an empty graph  $\mathcal{G}$  on vertex set  $\mathcal{V} = \{X_1, \dots, X_p\}$ 
2: Let  $S = \{1, \dots, p\}$  and  $\mathcal{T} = []$ 
3: repeat
4:    $H = []$ 
5:   for  $j \in S$  do
6:      $\hat{f}_j$ : Regress  $X^j$  on  $\{X_i\}_{i \in S \setminus \{j\}}$ 
7:      $\tilde{\zeta}_{\cdot j} = X_j - \hat{f}_j(X_i)$ 
8:      $h = \hat{\lambda}(\{X_i\}_{i \in S \setminus \{j\}}, \tilde{\zeta}_{\cdot j})$ 
9:      $H = [H, h]$ 
10:  end for
11:   $i^* = \arg \min_{i \in S} H$ 
12:   $S = S \setminus \{i^*\}$ 
13:   $\mathcal{T} = [i^*, \mathcal{T}]$ 
14: until  $|S| = 0$ 
15: for  $j \in \{2, \dots, p\}$  do
16:   for  $i \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\}$  do
17:     $\hat{f}_j$ : Regress  $X^j$  on  $\{X_k\}_{k \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\} \setminus \{i\}}$ 
18:     $\tilde{\zeta}_{\cdot j} = X_j - \hat{f}_j(X_i)$ 
19:    if  $\{X_k\}_{k \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\} \setminus \{i\}} \not\perp_P \tilde{\zeta}_{\cdot j}$  then
20:      Add  $X_j \rightarrow X_i$  to  $\mathcal{G}$ 
21:    end if
22:   end for
23: end for
24: Return  $\mathcal{G}$ 
```


ANM in action (1/4)

- ▶ Suppose the true graph on right;
- ▶ Assumptions: CMC, minimality, causal sufficiency.



ANM in action (2/4)

- ▶ Estimate $A, B, D \mapsto E$ and $\hat{\zeta}_e$
 - ▶ $H_1 = \hat{l}(\{A, B, D\}, \hat{\zeta}_e)$
- ▶ Estimate $A, D, E \mapsto B$ and $\hat{\zeta}_b$
 - ▶ $H_3 = \hat{l}(\{A, D, E\}, \hat{\zeta}_b)$
- ▶ Estimate $A, B, E \mapsto D$ and $\hat{\zeta}_d$
 - ▶ $H_2 = \hat{l}(\{A, B, E\}, \hat{\zeta}_d)$
- ▶ Estimate $B, D, E \mapsto A$ and $\hat{\zeta}_a$
 - ▶ $H_4 = \hat{l}(\{B, D, E\}, \hat{\zeta}_a)$

ANM in action (2/4)

- ▶ Estimate $A, B, D \mapsto E$ and $\hat{\zeta}_e$
 - ▶ $H_1 = \hat{l}(\{A, B, D\}, \hat{\zeta}_e)$
- ▶ Estimate $A, D, E \mapsto B$ and $\hat{\zeta}_b$
 - ▶ $H_3 = \hat{l}(\{A, D, E\}, \hat{\zeta}_b)$
- ▶ Estimate $A, B, E \mapsto D$ and $\hat{\zeta}_d$
 - ▶ $H_2 = \hat{l}(\{A, B, E\}, \hat{\zeta}_d)$
- ▶ Estimate $B, D, E \mapsto A$ and $\hat{\zeta}_a$
 - ▶ $H_4 = \hat{l}(\{B, D, E\}, \hat{\zeta}_a)$

$$4 = \mathit{Argmin}(H)$$
$$\mathcal{T} = [A]$$

ANM in action (3/4)

- ▶ Estimate $B, D \mapsto E$ and $\hat{\zeta}_e$
 - ▶ $H_1 = \hat{l}(\{B, D\}, \hat{\zeta}_e)$
- ▶ Estimate $D, E \mapsto B$ and $\hat{\zeta}_b$
 - ▶ $H_3 = \hat{l}(\{D, E\}, \hat{\zeta}_b)$
- ▶ Estimate $B, E \mapsto D$ and $\hat{\zeta}_d$
 - ▶ $H_2 = \hat{l}(\{B, E\}, \hat{\zeta}_d)$

ANM in action (3/4)

- ▶ Estimate $B, D \mapsto E$ and $\hat{\zeta}_e$
 - ▶ $H_1 = \hat{l}(\{B, D\}, \hat{\zeta}_e)$
- ▶ Estimate $D, E \mapsto B$ and $\hat{\zeta}_b$
 - ▶ $H_3 = \hat{l}(\{D, E\}, \hat{\zeta}_b)$
 - $\mathbf{1} = \text{Argmin}(H)$
 - $\mathcal{T} = [E, A]$
- ▶ Estimate $B, E \mapsto D$ and $\hat{\zeta}_d$
 - ▶ $H_2 = \hat{l}(\{B, E\}, \hat{\zeta}_d)$

ANM in action (3/4)

- ▶ Estimate $B, D \mapsto E$ and $\hat{\xi}_e$
 - ▶ $H_1 = \hat{l}(\{B, D\}, \hat{\xi}_e)$
- ▶ Estimate $D, E \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_3 = \hat{l}(\{D, E\}, \hat{\xi}_b)$
 - $1 = \text{Argmin}(H)$
 - $\mathcal{T} = [E, A]$
- ▶ Estimate $D \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_1 = \hat{l}(D, \hat{\xi}_b)$
- ▶ Estimate $B, E \mapsto D$ and $\hat{\xi}_d$
 - ▶ $H_2 = \hat{l}(\{B, E\}, \hat{\xi}_d)$
- ▶ Estimate $B \mapsto D$ and $\hat{\xi}_d$
 - ▶ $H_2 = \hat{l}(B, \hat{\xi}_d)$

ANM in action (3/4)

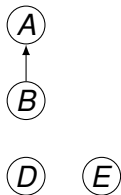
- ▶ Estimate $B, D \mapsto E$ and $\hat{\xi}_e$
 - ▶ $H_1 = \hat{l}(\{B, D\}, \hat{\xi}_e)$
- ▶ Estimate $D, E \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_3 = \hat{l}(\{D, E\}, \hat{\xi}_b)$
 $1 = \text{Argmin}(H)$
 $\mathcal{T} = [E, A]$
- ▶ Estimate $D \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_1 = \hat{l}(D, \hat{\xi}_b)$
 $2 = \text{Argmin}(H)$
 $\mathcal{T} = [D, E, A]$
- ▶ Estimate $B, E \mapsto D$ and $\hat{\xi}_d$
 - ▶ $H_2 = \hat{l}(\{B, E\}, \hat{\xi}_d)$
- ▶ Estimate $B \mapsto D$ and $\hat{\xi}_d$
 - ▶ $H_2 = \hat{l}(B, \hat{\xi}_d)$

ANM in action (3/4)

- ▶ Estimate $B, D \mapsto E$ and $\hat{\xi}_e$
 - ▶ $H_1 = \hat{l}(\{B, D\}, \hat{\xi}_e)$
 - ▶ Estimate $D, E \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_3 = \hat{l}(\{D, E\}, \hat{\xi}_b)$
 $1 = \text{Argmin}(H)$
 $\mathcal{T} = [E, A]$
 - ▶ Estimate $D \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_1 = \hat{l}(D, \hat{\xi}_b)$
 $2 = \text{Argmin}(H)$
 $\mathcal{T} = [D, E, A]$
 - ▶ Estimate $B, E \mapsto D$ and $\hat{\xi}_d$
 - ▶ $H_2 = \hat{l}(\{B, E\}, \hat{\xi}_d)$
 - ▶ Estimate $B \mapsto D$ and $\hat{\xi}_d$
 - ▶ $H_2 = \hat{l}(B, \hat{\xi}_d)$
- $\mathcal{T} = [B, D, E, A]$

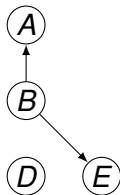
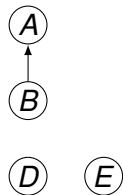
ANM in action (4/4)

$$\mathcal{T} = [B, D, E, A]$$



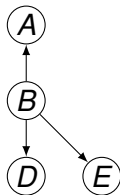
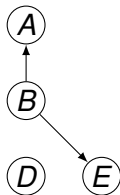
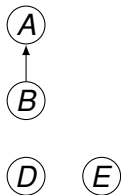
ANM in action (4/4)

$$\mathcal{T} = [B, D, E, A]$$



ANM in action (4/4)

$$\mathcal{T} = [B, D, E, A]$$



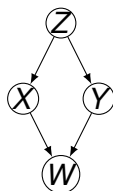
Exercise 1

After applying LiNGAM, how can you know if causal sufficiency is not respected?

Exercise 2

- ▶ Suppose the true graph on right;
- ▶ Assumptions: CMC, causal sufficiency, minimality;
- ▶ Generative process:

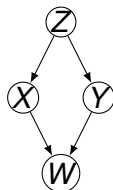
$$\begin{aligned}Z &= \zeta_z & \zeta_z &\sim U(0, 1); \\X &= a * Z + \zeta_x & \zeta_x &\sim U(0, 1); \\Y &= b * Z + \zeta_y & \zeta_y &\sim U(0, 1); \\W &= c * X - d * Y + \zeta_w & \zeta_w &\sim N(0, 1).\end{aligned}$$



- ▶ Given a compatible distribution what would be the output of the LiNGAM algorithm? And what about the ANM algorithm?

Exercise 3

- ▶ Suppose the true graph on right;
- ▶ Assumptions: CMC, causal sufficiency, minimality;
- ▶ Generative process:



$$\begin{aligned} Z &= \zeta_z & \zeta_z &\sim U(0, 1); \\ X &= Z^2 + \zeta_x & \zeta_x &\sim U(0, 1); \\ Y &= Z^3 + \zeta_y & \zeta_y &\sim U(0, 1); \\ W &= XY + \zeta_w & \zeta_w &\sim U(0, 1). \end{aligned}$$

- ▶ Given a compatible distribution what would be the output of the LiNGAM algorithm? And what about the ANM algorithm?

Table of content

Preliminaries

Bivariate causal discovery

Multivariate causal discovery

Conclusion

Conclusion

- ▶ Under linear non gaussian models noise-based methods can discover the causal graph.
- ▶ Under non-linear additive noise models noise-based methods can discover the causal graph.
- ▶ Advantages:
 - ▶ Can discovery the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.

Conclusion

- ▶ Under linear non gaussian models noise-based methods can discover the causal graph.
- ▶ Under non-linear additive noise models noise-based methods can discover the causal graph.
- ▶ Advantages:
 - ▶ Can discovery the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.

Conclusion

- ▶ Under linear non gaussian models noise-based methods can discover the causal graph.
- ▶ Under non-linear additive noise models noise-based methods can discover the causal graph.
- ▶ Advantages:
 - ▶ Can discovery the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.

Conclusion

- ▶ Under linear non gaussian models noise-based methods can discover the causal graph.
- ▶ Under non-linear additive noise models noise-based methods can discover the causal graph.
- ▶ Advantages:
 - ▶ Can discovery the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.

Some extensions

- ▶ Without causal sufficiency if linear relations;
- ▶ Extension to discrete additive noise models;
- ▶ Post non linear relations;
- ▶ Time series.

Some extensions

- ▶ Without causal sufficiency if linear relations;
- ▶ Extension to discrete additive noise models;
- ▶ Post non linear relations;
- ▶ Time series.

Some extensions

- ▶ Without causal sufficiency if linear relations;
- ▶ Extension to discrete additive noise models;
- ▶ Post non linear relations;
- ▶ Time series.

Some extensions

- ▶ Without causal sufficiency if linear relations;
- ▶ Extension to discrete additive noise models;
- ▶ Post non linear relations;
- ▶ Time series.

Direct inspirations

1. *Elements of causal inference*, J. Peters, D. Janzing , B. Schölkopf. MIT Press, 2nd edition, 2017
2. *DirectLiNGAM: A Direct Method for Learning a Linear Non-Gaussian Structural Equation Model*, S. Shimazu, T. Inazumi, Y. Sogawa, A. Hyvarinen, Y. Kawahara, T. Washio, P. Hoyer, K. Bollen. JMLR, 2011
3. *Nonlinear causal discovery with additive noise models*, P. Hoyer, D. Janzing, J. Mooij, J. Peters, B. Schölkopf. Neurips, 2008
4. *Causal Discovery with Continuous Additive Noise Models*, J. Peters, J. Mooij, D. Janzing, B. Schölkopf. JMLR, 2014

References (2/2)

Additional readings

1. *Causal inference from noise*, N. Climenhaga, L. DesAutels, G. Ramsey. Noûs, 2019
2. *On the logic of causal models*, D. Geiger, J. Pearl. In Proceedings of the Fourth Annual Conference on Uncertainty in Artificial Intelligence, 1990
3. *A Linear Non-Gaussian Acyclic Model for Causal Discovery*, S. Shimazu, P. Hoyer, A. Hyvarinen, A. Kerminen. JMLR, 2006
4. *Analyse générale des liaisons stochastiques.*, G. Darmais. Review of the International Statistical Institute, 1953
5. *On a property of the normal distribution*, W. P. Skitovitch. Doklady Akademii Nauk SSSR, 89:217–219, 1953
6. *Causal Inference on Time Series using Restricted Structural Equation Models*, J. Peters, D. Janzing, B. Schölkopf. Neurips, 2013