

Causal discovery: noise-based methods

Charles K. Assaad, Emilie Devijver, Eric Gaussier

emilie.devijver@univ-grenoble-alpes.fr

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

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, 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

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)

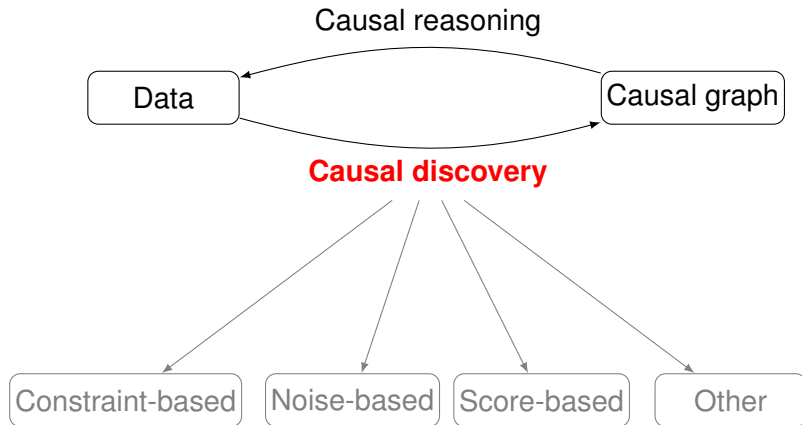
Independent Mechanism Principle

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

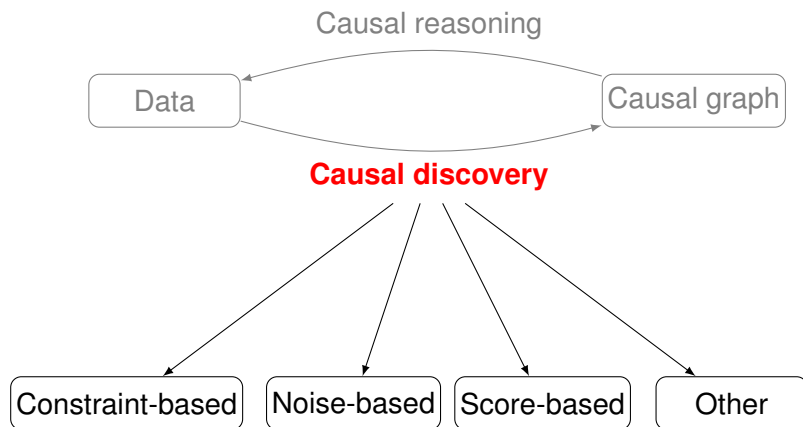
- ▶ Independence of noises, conditional independence of structures
- ▶ Independence of information contained in mechanisms

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 .

Causal discovery



Causal discovery



Recap: Constraint-based

- ▶ Independence based algorithm to infer the causal graph
- ▶ Need to test independence → need of large sample size
- ▶ Faithfulness
- ▶ Identify the Markov equivalence class

Can we do better than identify the Markov equivalence class?

Not without parametric assumptions!

Recap: Constraint-based

- ▶ Independence based algorithm to infer the causal graph
- ▶ Need to test independence → need of large sample size
- ▶ Faithfulness
- ▶ Identify the Markov equivalence class

Can we do better than identify the Markov equivalence class?

Not without parametric assumptions!

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

Example with two nodes

Consider two variables X and Y , and assume we have infinite data and can deduce the joint distribution.

$$X \rightarrow Y \quad \text{or} \quad Y \rightarrow X$$

Markov equivalence class: $X - Y$

Example with two nodes

Consider two variables X and Y , and assume we have infinite data and can deduce the joint distribution.

$$X \rightarrow Y \quad \text{or} \quad Y \rightarrow X$$

Markov equivalence class: $X - Y$

Non-identifiability of two-node graphs

Proposition (Peters et al., 2017, p44)

For every joint distribution $P(x, y)$ on two real-valued random variables, there is an SCM in either direction that generates data consistent with $P(x, y)$.

Mathematically, there exists a function f_Y such that

$$Y = f_Y(X, U_Y), X \perp\!\!\!\perp U_Y$$

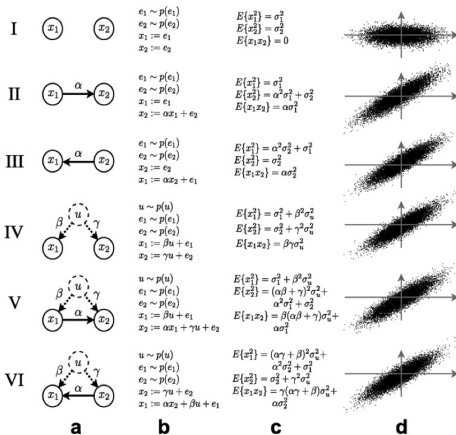
and there exists a function f_X such that

$$X = f_X(Y, U_X), Y \perp\!\!\!\perp U_X$$

where U_Y and U_X are real-valued random variables.

Non-identifiability

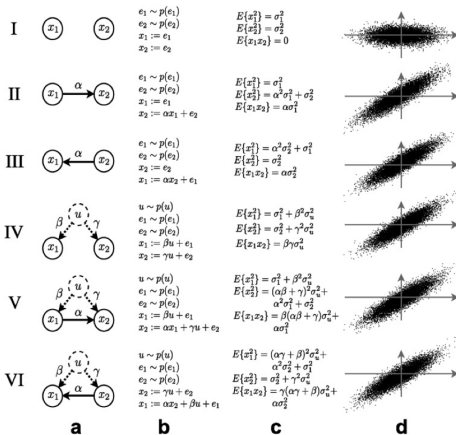
Ex from Hoyer et al. (2008)



The Markov equivalence class is the best we can do!

Non-identifiability

Ex from Hoyer et al. (2008)



The Markov equivalence class is the best we can do!

Non-identifiability

Multinomial or Gaussian distribution

Theorem (Markov Completeness, Meek (1995), Geiger and Pearl (1988))

If we have multinomial distributions or linear Gaussian structural equations, we can only identify a graph up to its Markov equivalence class.

The intuition behind the noise

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

The intuition behind the noise

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

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

The intuition behind the noise

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

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

$$Y := 2X + \zeta_y \quad \text{or} \quad X := \frac{Y}{2} + \zeta_x?$$

The intuition behind the noise

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

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

$$Y := 2X + \zeta_y \quad \text{or} \quad X := \frac{Y}{2} + \zeta_x?$$

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

The intuition behind the noise

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

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

$$Y := 2X + \zeta_y \quad \text{or} \quad X := \frac{Y}{2} + \zeta_x?$$

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

X	Y	$\zeta_y = Y - 2X$	$\zeta_x = X - Y/2$
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	$1 \in \{-1, 0, 1\}$	$-0.5 \notin \{-1, 0, 1\}$

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

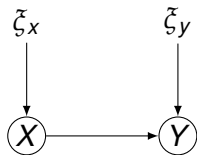
The post non linear case

In practice

Multivariate case

Conclusion

The linear case



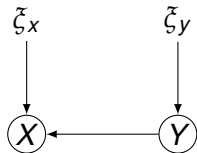
$$M_1 : \begin{cases} X := \zeta_x \\ Y := aX + \zeta_y \end{cases}$$

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

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

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

$$Y = aX + \zeta_y$$
$$\zeta_x = (1 - ba)X - b\zeta_y$$

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

Theorem (Darmois-Skitovich, 1953, 1954)

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 case

$$Y = aX + \zeta_y$$
$$\zeta_x = (1 - ba)X - b\zeta_y$$

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

Theorem (Darmois-Skitovich, 1953, 1954)

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

Theorem (Identifiability of linear non-Gaussian models)

Assume that $P(X, Y)$ admits the linear model

$$Y := aX + \zeta_y, \quad X \perp\!\!\!\perp_P \zeta_y,$$

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

$$X := bY + \zeta_x, \quad Y \perp\!\!\!\perp_P \zeta_x,$$

if and only if ζ_y and X are Gaussian.

It can be extended to multiple variables (see Shimizu et al. (2006), Peters et al. (2017)).

The linear non gaussian case

Theorem (Identifiability of linear non-Gaussian models)

Assume that $P(X, Y)$ admits the linear model

$$Y := aX + \zeta_y, \quad X \perp\!\!\!\perp_P \zeta_y,$$

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

$$X := bY + \zeta_x, \quad Y \perp\!\!\!\perp_P \zeta_x,$$

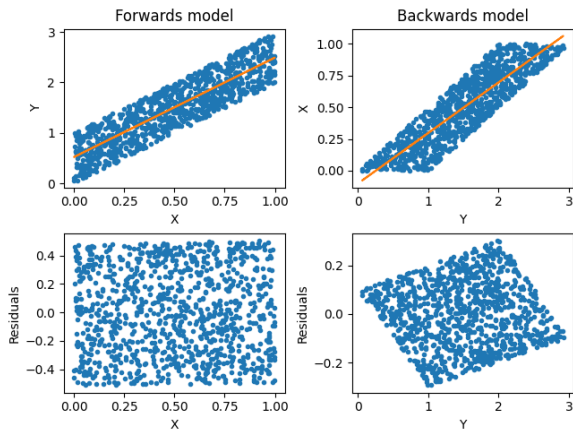
if and only if ζ_y and X are Gaussian.

It can be extended to multiple variables (see Shimizu et al. (2006), Peters et al. (2017)).

The linear non gaussian case

Example:

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



The linear non gaussian case

Ex from Hoyer et al. (2008)

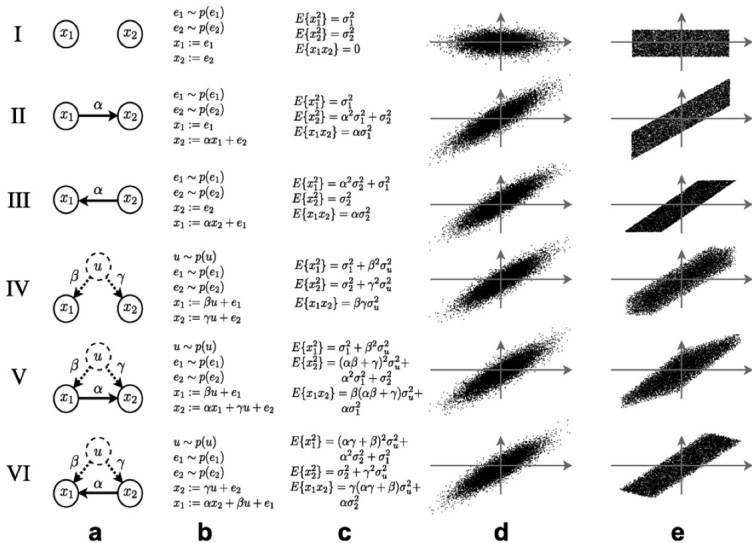


Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

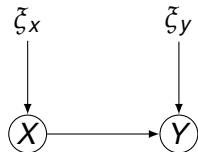
In practice

Multivariate case

Conclusion

The non linear ANM case

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 ANM case

Theorem (Identifiability of additive noise models, Hoyer et al, 2008)

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 f_X 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.

The non linear ANM case

Theorem (Identifiability of additive noise models, Hoyer et al, 2008)

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 f_x 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.

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

- ▶ 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.

Be careful! the ANM is not closed under marginalization: with latent variables, if you assume ANM over $X \cup L$, then the marginal model over X may no longer be in the ANM class.

The non linear case

- ▶ 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.

Be careful! the ANM is not closed under marginalization: with latent variables, if you assume ANM over $X \cup L$, then the marginal model over X may no longer be in the ANM class.

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

The post non linear case

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

Zhang and Hyvarinen (2009) have provided conditions under which this causal model is identifiable.

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

Causal order discovery procedure in the bivariate case

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

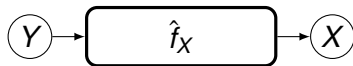
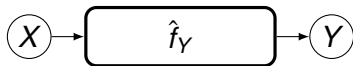
Procedure:

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 :

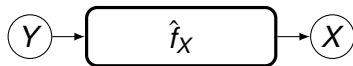
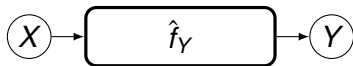


Causal order discovery procedure in the bivariate case

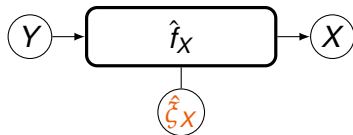
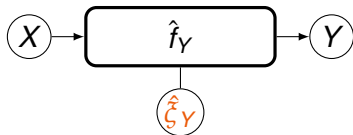
Given $P(X, Y)$ and a dependence estimator $\hat{\gamma}$

Procedure:

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



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

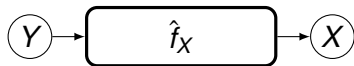
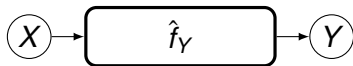


Causal order discovery procedure in the bivariate case

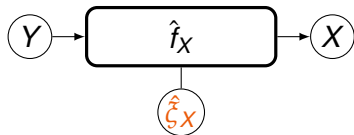
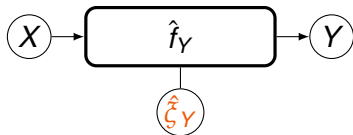
Given $P(X, Y)$ and a dependence estimator $\hat{\lambda}$

Procedure:

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



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



3. Order:

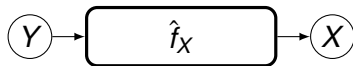
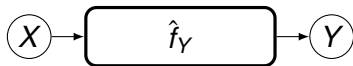
- ▶ $\mathcal{T} = [X, Y]$ if $\hat{\lambda}(X, \hat{\zeta}_Y) < \hat{\lambda}(Y, \hat{\zeta}_X)$
- ▶ $\mathcal{T} = [Y, X]$ if $\hat{\lambda}(Y, \hat{\zeta}_X) < \hat{\lambda}(X, \hat{\zeta}_Y)$

Causal order discovery procedure in the bivariate case

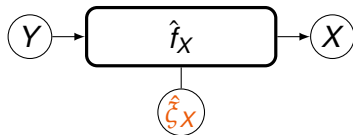
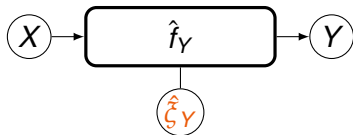
Given $P(X, Y)$ and a dependence estimator $\hat{\lambda}$

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

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

Assumptions

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 $i < j$.

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} .

The graph does not contain dependencies not present in the observational data.

Assumptions

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 $i < j$.

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} .
The graph does not contain dependencies not present in the observational data.

Assumptions

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 $i < j$.

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} .
The graph does not contain dependencies not present in the observational data.

Remark: under Markov condition, faithfulness \implies minimality.

The algorithms

1. Topological ordering
2. Pruning

The algorithms

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: $\hat{\xi}_{ij} = X_j - \frac{\text{cov}(X_i, X_j)}{\text{var}(X_i)} X_i$
 - 8: $h = \sum_{j \in S \setminus \{i\}} \hat{\lambda}(X_i, \hat{\xi}_{ij})$
 - 9: $H = [H, h]$
 - 10: $i^* = \arg \min_{i \in S} H$
 - 11: $S = S \setminus \{i^*\}$
 - 12: $\mathcal{T} = [\mathcal{T}, i^*]$
 - 13: $\forall j \in S, X_j = \hat{\xi}_{j^*}$
 - 14: **until** $|S| = 0$
 - 15: Append (\mathcal{T}, S_0)
 - 16: 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.
 - 17: **if** $a_{i,j} > 0$ **then**
 - 18: Add $X_i \rightarrow X_j$ to \mathcal{G}
 - 19: **Return** \mathcal{G}
-

Algorithm 2 ANM

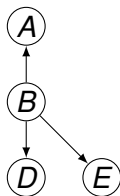
Input: $P(\mathcal{V})$

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: $\hat{\xi}_{\cdot j} = X_j - \hat{f}_j(X_i)$
 - 8: $h = \hat{\lambda}(\{X_i\}_{i \in S \setminus \{j\}}, \hat{\xi}_{\cdot j})$
 - 9: $H = [H, h]$
 - 10: $i^* = \arg \min_{i \in S} H$
 - 11: $S = S \setminus \{i^*\}$
 - 12: $\mathcal{T} = [\mathcal{T}, i^*]$
 - 13: **until** $|S| = 0$
 - 14: **for** $j \in \{2, \dots, p\}$ **do**
 - 15: **for** $i \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\}$ **do**
 - 16: \hat{f}_j : Regress X^j on $\{X_k\}_{k \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\} \setminus \{i\}}$
 - 17: $\hat{\xi}_{\cdot j} = X_j - \hat{f}_j(X_i)$
 - 18: **if** $\{X_k\}_{k \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\} \setminus \{i\}} \not\perp_P \hat{\xi}_{\cdot j}$ **then**
 - 19: Add $X_i \rightarrow X_j$ to \mathcal{G}
 - 20: **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)

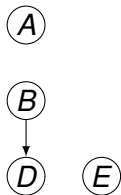
- ▶ Estimate $B, D \mapsto E$ and $\hat{\xi}_e$
 - ▶ $H_1 = \hat{l}(\{B, D\}, \hat{\xi}_e)$
- ▶ Estimate $D, E \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_3 = \hat{l}(\{D, E\}, \hat{\xi}_b)$
 $1 = \text{Argmin}(H)$
 $\mathcal{T} = [E, A]$
- ▶ Estimate $D \mapsto B$ and $\hat{\xi}_b$
 - ▶ $H_1 = \hat{l}(D, \hat{\xi}_b)$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

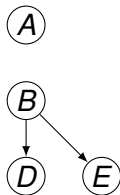
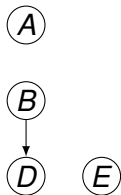
ANM in action (4/4)

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



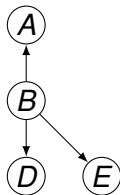
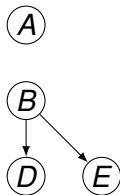
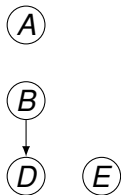
ANM in action (4/4)

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



ANM in action (4/4)

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



Exercise 1

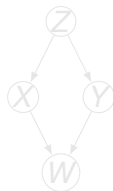
- ▶ Suppose the true graph on right;
- ▶ Assumptions: CMC, causal sufficiency, minimality;
- ▶ Generative process:

$$Z = \zeta_Z \quad \zeta_Z \sim U(0, 1);$$

$$X = a * Z + \zeta_X \quad \zeta_X \sim U(0, 1);$$

$$Y = b * Z + \zeta_Y \quad \zeta_Y \sim U(0, 1);$$

$$W = c * X - d * Y + \zeta_W \quad \zeta_W \sim N(0, 1).$$



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

Exercise 1

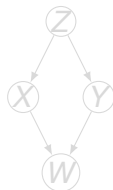
- ▶ Suppose the true graph on right;
- ▶ Assumptions: CMC, causal sufficiency, minimality;
- ▶ Generative process:

$$Z = \zeta_Z \quad \zeta_Z \sim U(0, 1);$$

$$X = a * Z + \zeta_X \quad \zeta_X \sim U(0, 1);$$

$$Y = b * Z + \zeta_Y \quad \zeta_Y \sim U(0, 1);$$

$$W = c * X - d * Y + \zeta_W \quad \zeta_W \sim N(0, 1).$$



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

And what about PC?

Exercise 1

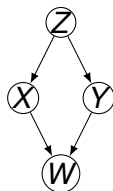
- ▶ Suppose the true graph on right;
- ▶ Assumptions: CMC, causal sufficiency, minimality;
- ▶ Generative process:

$$Z = \zeta_Z \quad \zeta_Z \sim U(0, 1);$$

$$X = a * Z + \zeta_X \quad \zeta_X \sim U(0, 1);$$

$$Y = b * Z + \zeta_Y \quad \zeta_Y \sim U(0, 1);$$

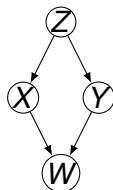
$$W = c * X - d * Y + \zeta_W \quad \zeta_W \sim N(0, 1).$$



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

Exercise 2

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



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

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

Exercise 3

Why is faithfulness needed for constraint-based methods whereas noise-based methods only need minimality?

Table of content

Reminder

Problem statement

The linear case

The non linear ANM case

The post non linear case

In practice

Multivariate case

Conclusion

Conclusion

- ▶ Noise-based methods can discover the causal graph
 - ▶ Under linear non gaussian models
 - ▶ Under non-linear additive noise models
- ▶ Advantages:
 - ▶ Can discover the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.
- ▶ Extensions
 - ▶ Without causal sufficiency if linear relations;
 - ▶ With cyclic graphs;
 - ▶ Extension to discrete additive noise models;
 - ▶ Time series.

Conclusion

- ▶ Noise-based methods can discover the causal graph
 - ▶ Under linear non gaussian models
 - ▶ Under non-linear additive noise models
- ▶ Advantages:
 - ▶ Can discover the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.
- ▶ Extensions
 - ▶ Without causal sufficiency if linear relations;
 - ▶ With cyclic graphs;
 - ▶ Extension to discrete additive noise models;
 - ▶ Time series.

Conclusion

- ▶ Noise-based methods can discover the causal graph
 - ▶ Under linear non gaussian models
 - ▶ Under non-linear additive noise models
- ▶ Advantages:
 - ▶ Can discover the true graph;
 - ▶ Faithfulness is not needed.
- ▶ Drawbacks:
 - ▶ Semi parametric assumptions;
 - ▶ Need large sample size.
- ▶ Extensions
 - ▶ Without causal sufficiency if linear relations;
 - ▶ With cyclic graphs;
 - ▶ Extension to discrete additive noise models;
 - ▶ Time series.

References (1/2)

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. Shimizu, 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