# **Outline**

### ☑ 背景知识

### ☑ 因果发现

- ・ 什么是因果发现
- · 经典方法:基于约束的方法、基于因果函数的方法
- ・ 研究进展: 隐变量问题、非独立同分布问题
- ・ 应用探索: 故障检测

### ☑ 因果性学习







☑ 因果发现:回答"为什么?"











☑ **基于实验的方法**:干预原因,结果会发生改变



# 如何发现因果关系?基于观察数据的方法



☑ 基于观察数据的方法: 观察数据+因果假设⇒因果模型







5



- ☑ 基于约束的因果发现方法
- ☑ 基于函数的因果发现方法
- ☑ 混合型因果发现方法

# **Structural Causal Models and Graphical Causal Models**

- DMIR
- ☑ The structural causal model (SCM) is a framework that can be used for multivariate analysis, which can be used to describe the real-world related variables and their interactions.

$$\boxdot X_i = f_i(pa_i, E_i), i = 1, 2, \dots, n$$

- $pa_i$  : parents of  $X_i$
- $E_i$ : exogenous variables / errors / disturbances
- Each equation represents an autonomous mechanism
- Describes how nature assigns values to variables of interest

#### SCM

$$X = \{X_1, X_2, X_3, X_4\}, E = \{E_1, E_2, E_3, E_4\}, F = \{f_3, f_4\}$$
$$X_1 = E_1$$
$$X_2 = E_2$$

$$X_3 = f_3(X_1, X_2, E_3)$$
  
 $X_4 = f_4(X_3, E_4)$ 

Graphical Causal Model: Directed Acyclic Graphs (DAG)

Nodes: 
$$\{X_1, X_2, X_3, X_4\}$$
 $X_1$  $X_2$ Edges:  $\{X_1 \rightarrow X_3,$  $X_3$  $X_2 \rightarrow X_3,$  $X_3$  $X_3 \rightarrow X_4\}$  $X_4$ 



✓ Causal Markov Assumption: A variable X is independent of every other variable (except X's effects) conditional on all of its direct causes.



$$\Leftrightarrow \quad x_4 \amalg \{x_3, x_6, x_5\} \mid \{x_1, x_2\}$$

**Causal Faithfulness Assumption**: for all observed variables,  $X_i$  is independent of  $X_j$  conditional on variables **Z** if and only if the Markov Assumption for G entails such conditional independencies.

### **Constraint-based methods**



☑ (Conditional) Independence Test





#### ☑ PC (Peter-Clark from CMU)

<i>X</i> <sub>1</sub>	$X_2$	$X_3$	$X_4$	$X_5$
-1.1	1	1.3	0.2	-0.7
2.1	2	3.1	-1.3	-1.6
3.1 2.3	4.2 -0.6	-2.6 -3.5	0.6 0.8	2.1
1.3	-1.7	0.9	2.4	-1.4
-1.8	0.9	-1.3	0.9	0.7
(a) 数据				

Spirtes, P., Glymour, C. N., Scheines, R., & Heckerman, D. Causation, prediction, and search. MIT press, 2000.

### **Constraint-based methods: PC**







 $\blacksquare$  How to deal with latent confounders?

• if there is a latent confounder L behind  $X_3$  and  $X_4$ 



Spirtes, P., Glymour, C. N., Scheines, R., & Heckerman, D. Causation, prediction, and search. MIT press, 2000.



☑ FCI (Fast Causal Inference): allows Confounders

• Results represented by PAGs (Partial Ancestral Graphs)

$$X_1$$
  $X_2$   $X_1$  and  $X_2$  are not adjacent

$$X_1 \longrightarrow X_2$$
  $X_2$  is not an ancestor of  $X_1$ 

$$X_1 \circ \cdots \circ X_2$$
 No set d-separates  $X_2$  and  $X_1$ 

$$X_1 \longrightarrow X_2 \qquad X_1 \text{ is a cause of } X_2$$

 $X_1 \blacktriangleleft$ 

 $\longrightarrow$   $X_2$  There is a latent common cause of  $X_1$  and  $X_2$ 



Spirtes, P., Glymour, C. N., Scheines, R., & Heckerman, D. Causation, prediction, and search. MIT press, 2000.



☑ Limitations of constraint-based methods



Markov Equivalence Class

Problem: Cannot identify the structures belonging to the Markov Equivalence Class

Spirtes, P., Glymour, C. N., Scheines, R., & Heckerman, D. Causation, prediction, and search. MIT press, 2000.



- ☑ 基于约束的因果发现方法
- ☑ 基于函数的因果发现方法
- ☑ 混合型因果发现方法

## **Can we directly distinguish cause from effect?**





 $X_1 \rightarrow X_2 \text{ or } X_2 \rightarrow X_1?$ 



### ☑ Considering the data generating process, effect generated from causes and noises,

represented with functional causal model:

$$Y = f(X, E)$$

### ☑ Introducing additional assumptions

• Independent noise assumption: Independence between the causes X and noises E

• Independent mechanism assumption: independence between the causes X and process f









 $\ensuremath{\boxtimes}$  Causal Asymmetry in the Linear non-Gaussian Case

• Data generated by Y = aX + E (*i.e.*,  $X \to Y$ )

 $\square$  (**X**, *Y*) follows the IN condition iff regression residual *Y* –  $\tilde{\omega}^T \mathbf{X}$  is independent from **X** 



# **Causal Function based method: LiNGAM**

DMIR

☑ Under the above assumptions, the LiNGAM can be expressed as

#### $\mathbf{X} = \mathbf{B}\mathbf{X} + \mathbf{E}$

- *X* is a p-dimensional random vector, representing the observed variable.
- **B** is  $p \times p$ -dimensional matrix, which represents the connection weight between the observed variables.
- *E* is a p-dimensional non-Gaussian random noise variable.
- $\square$  Because of the DAG assumption, there exists a permutation matrix  $P \in \mathbb{R}^{m \times m}$  such that **B**'

 $= PBP^{T}$  is a strict lower triangular matrix and diagonal elements are all 0

# **LiNGAM: Independent Component Analysis**





#### Assumptions in ICA

- At most one of  $S_i$  is Gaussian
- Size(X) >= Size(S), and **A** is of full column rank

Hyvärinen et al., Independent Component Analysis, 2001

Then A can be estimated up to column scale and permutation indeterminacies



#### ☑ LINGAM:

 $\mathbf{X} = \mathbf{B}\mathbf{X} + \mathbf{E}$ 

☑ ICA:

Y = WXB = I - W

#### ☑ An example

$$\begin{bmatrix} E_1 \\ E_3 \\ E_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -0.5 & 0 \\ 0.2 & 0 \end{bmatrix} \cdot \begin{bmatrix} X_2 \\ X_3 \\ X_1 \end{bmatrix}$$
$$\Leftrightarrow \begin{cases} X_2 = E_1 \\ X_3 = 0.5X_2 + E_3 \\ X_1 = -0.2X_2 + 0.3X_3 + E_2 \end{cases}$$

So we have the causal relation:



Shimizu S, Hoyer PO, Hyvärinen A, et al. A linear non-Gaussian acyclic model for causal discovery[J]. Journal of Machine Learning Research, 2006, 7(10).

### **Causal Function based method: ANM**

☑ Hoyer et al. proved that nonlinear functions can play a similar role to non-Gaussian models, which can be used to identify causal directions.

 $Y = f(X) + E \quad with$  $E \bot\!\!\!\bot X$ 



(Hoyer et al., 2009)





#### ☑ Motivation: Non-Transitivity of Nonlinear Causal Model

• Let the direct cause in  $X_1 \rightarrow X_2 \rightarrow X_3$  satisfy the additive noise model (ANM):

 $\begin{cases} X_1 \sim U(-0.5, 0.5) \\ X_2 = 2 \tanh(5X_1) + N_2 \\ X_3 = \left(\frac{X_2}{2}\right)^3 + N_3 \end{cases}$ 

• However, the causal influence  $X_1 \rightarrow X_3$  does not necessarily follow the additive noise model,



and we might not identify the causal by simply test the independence between cause and noise.

### Causal Function based method: cascade Additive Noise Model (CANM)

- DMI
- ☑ There exists a sequence of unmeasured intermediate variables Z between X and Y, where  $N_1$ ,  $N_2$ ,  $N_3$ ,  $\epsilon$  are the additive noise at each direct cause.



Figure: Illustration of the CANM

☑ Formally:

$$\begin{cases} Z_1 = f_1(X) + N_1 \\ Z_t = f_t(\mathbf{Z}_{pa(t)}) + N_t \\ Y = f_{T+1}(\mathbf{Z}_{pa(y)}) + \epsilon \end{cases}$$

#### $\blacksquare$ How to solve it?



 $\square$  Given data set  $\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^{m}$ , the marginal likelihood is given as follow:

$$\log \prod_{i=1}^{m} \int p_{\theta} \left( x^{(i)}, y^{(i)}, z \right) dz \Leftrightarrow \log \prod_{i=1}^{m} \int p_{\theta} \left( x^{(i)}, \epsilon^{(i)}, n \right) dn$$

① Decomposes the joint likelihood based on the Markov condition

② Applies the independence property between the cause and the noise, i.e.,

$$p(Z_t | \mathbf{Z}_{pa(t)}) = p(N_t = Z_t - f_t(\mathbf{Z}_{pa(t)}) | \mathbf{Z}_{pa(t)}) \underline{\mathbf{Z}_{pa(t)} \perp N_t} p(N_t = Z_t - f_t(\mathbf{Z}_{pa(t)}))$$

③ At the same time, we replace  $d\mathbf{z} = d\mathbf{n}$  and rewrite function  $f_{T+1}(\mathbf{Z}_{pa(y)})$  as  $f(X, \mathbf{N})$ .

#### ☑ The unobserved intermediate variables z can be replaced by n.

# **CANM: variational solution**



 $\square$  Purpose: Find  $q_{\phi}(\mathbf{N}|X,Y)$  to approximate  $p_{\theta}(\mathbf{N}|X,Y)$ 

☑ Action: Optimize a variational lower bound (ELBO) of the marginal log-likelihood

$$\log \prod_{\substack{i=1 \\ m \ i=1}}^{m} \int p_{\theta} \left( x^{(i)}, \epsilon^{(i)}, n \right) dn$$
  

$$\geq \sum_{i=1}^{m} \log p \left( x^{(i)} \right) - KL(q_{\phi} \left( n | x^{(i)}, y^{(i)} \right) \| p_{\theta}(n) \right) + E_{n \sim q_{\phi}(n | x^{(i)}, y^{(i)})} \left[ \log p \left( \epsilon^{(i)} = y^{(i)} - f \left( x^{(i)}, n; \theta \right) \right) \right]$$

 $\square$  The lower bound is tight at the maximum of ELBO i.e.  $q_{\phi}(\mathbf{N}|X,Y)$  equal  $p_{\theta}(\mathbf{N}|X,Y)$ 

 $\square$  Steps of VAE to optimize  $\mathcal{L}(\theta, \phi; X, Y)$ 

- 1. Encode data into  $\mu$ ,  $\sigma$ .
- 2. Sample the *n* through  $\mu + \sigma \odot u$ , where  $u \sim \mathcal{N}(0, I)$
- 3. Reconstruct y using x, n.





 $Z_2$ 

 $Z_3$ 

 $Z_1$ 

Х





 $\square$  Model selection: select the best number of latent variables.

☑ To estimate the marginal likelihood to identify the causal relationship.





 $\square$  Theorem 1. Let  $X \rightarrow Y$  follow the cascade additive noise model, while there exists a backward model following the same form, i.e.

 $Y = f(X, N) + \epsilon,$  X, N, and  $\epsilon$  are independent,  $X = g(Y, \widehat{N}) + \hat{\epsilon},$  Y,  $\widehat{N}$ , and  $\hat{\epsilon}$  are independent,

then the noise distribution of the reverse direction  $p_{\hat{\epsilon}}$  must be

$$p_{\hat{\epsilon}}(\hat{\epsilon}) = \int e^{2\pi i \hat{\epsilon} \cdot \nu} \frac{\int \int p(x) p(n) p_{\epsilon}(y - f(x, n)) e^{-2\pi i x \cdot \nu} dn dx}{p(y) \int p(\hat{n}) e^{-2\pi i g(y, \hat{n}) \cdot \nu} d\hat{n}} d\nu_{\epsilon}$$

where f, g denote the function implied by the cascade process.

• The main result of our main theorem is that, If the model is non-identifiable, one strict condition must hold:

$$\begin{aligned} \forall y_1, y_2, \int e^{2\pi i\hat{\epsilon}\cdot\nu} \frac{\int \int p(x)p(\mathbf{n})p_{\epsilon}(y_1 - f(x, \mathbf{n}))e^{-2\pi ix\cdot\nu}d\mathbf{n}dx}{p(y_1)\int p(\hat{\mathbf{n}})e^{-2\pi ig(y_1, \hat{\mathbf{n}})\cdot\nu}d\hat{\mathbf{n}}}d\nu \\ &= \int e^{2\pi i\hat{\epsilon}\cdot\nu} \frac{\int \int p(x)p(\mathbf{n})p_{\epsilon}(y_2 - f(x, \mathbf{n}))e^{-2\pi ix\cdot\nu}d\mathbf{n}dx}{p(y_2)\int p(\hat{\mathbf{n}})e^{-2\pi ig(y_2, \hat{\mathbf{n}})\cdot\nu}d\hat{\mathbf{n}}}d\nu \end{aligned}$$

#### Intuitively, such a condition holds only in restrictive cases.

Cai R, Qiao J, Zhang K, et al. Causal discovery with cascade nonlinear additive noise models. IJCAI, 2019.

### **CANM: experiments - synthetic data**



- ☑ Depth: as the depth increases, the accuracy of CANM is stable and around 90% accuracy with a slight decrease, while the performance of the rest methods decreases rapidly as the depth grows.
- ☑ Sample Size: as the sample size grows the accuracy increase, and even in the small sample size, CANM still outperforms the other methods.
- ☑ Fixed Structure: the variance of the likelihood decreases and the accuracy increases as the sample size grows.







Figure 6: Sensitivity to Sample in a Fixed Structure.

# **CANM: experiments - real world data**



#### $\ensuremath{\boxtimes}$ The electricity consumption dataset

- Hour of day  $\rightarrow$  Temperature  $\rightarrow$  electricity load
  - This data might exist more than one unmeasured variables because in the same hour of day there have different electricity load and the reason could be the season.
  - CANM successfully catch this latent variable because the prediction of electricity load, the red point, separating into both upper and lower parts.

### ☑ Stock Market Dataset

- Hutchison  $\rightarrow$  Cheung Kong  $\rightarrow$  Sun Hung Kai.
  - The fitted intermediate variable from CANM has a high correction ( $\rho = 0.54$ ) with the ground truth (Cheung Kong)



Cai R, Qiao J, Zhang K, et al. Causal discovery with cascade nonlinear additive noise models. IJCAI, 2019.

### **Causal Function based method: PNL**

- ☑ LiNGAM algorithms can only solve linear problems. For non-linear problems, Zhang et al. proposed PNL, the post-NonLinear method.
- $\square$  In the PNL model, assuming that there is a causal relationship  $v_i \rightarrow v_j$ , it can be expressed as

 $v_j = f_2(f_1(v_i) + n_j)$  (1)

- $v_i$  and  $n_j$  are independent of each other
- $f_1$  is an unconstant smooth function
- $f_2$  is a reversible smooth function and  $f'_2 \neq 0$



## **PNL: All non-identifiable cases**



☑ Causal direction is generally identifiable if the data were generated according to

 $X_2 = f_2(f_1(X_1) + E).$ 





Table 1: All situations in which the PNL causal model is not identifiable.

	$p_{e_2}$	$p_{t_1} (t_1 = g_2^{-1}(x_1))$	$h = f_1 \circ g_2$	Remark	
Ι	Gaussian	Gaussian	linear	$h_1$ also linear	
II	log-mix-lin-exp	log-mix-lin-exp	linear	$h_1$ strictly monotonic, and $h'_1 \rightarrow$	
				0, as $z_2 \to +\infty$ or as $z_2 \to -\infty$	
III	log-mix-lin-exp	one-sided asymptoti-	h strictly monotonic,	—	
		cally exponential (but	and $h' \to 0$ , as $t_1 \to 0$		
		not log-mix-lin-exp)	$+\infty$ or as $t_1 \to -\infty$		
IV	log-mix-lin-exp	generalized mixture of	Same as above	—	
		two exponentials			
V	generalized mixture	two-sided asymptoti-	Same as above	—	
	of two exponentials	cally exponential 📐			

$$p_v \propto (c_1 e^{c_2 v} + c_3 e^{c_4 v})^{c_5}$$

Zhang K, Hyvarinen A. On the identifiability of the post-nonlinear causal model. UAI, 2009.

### **Causal Function based methods**



☑ If we add new assumptions to the PNL model, we will get the LiNGAM and ANM.



## **Causal Function based methods: in noiseless case**



- $\square$  Problem: how to infer whether Y = f(X) or  $X = f^{-1}(X)$  is the right causal model?
- ☑ An example of independent mechanism



☑ Asymmetry holds

- In the causal direction: independence between X and f(X) holds
- In the anti-causal direction: independence between Y and f(Y) does not hold



 $\boxdot If X \to Y then$ 

$$\int \log |f'(x)| \, p(x) dx \, \leq \int \log \left| f^{-1'}(y) \right| \, p(y) dy$$

☑ empirical estimator

$$\hat{C}_{X \to Y} \coloneqq \frac{1}{m} \sum_{j=1}^{m} \log \left| \frac{y_{j+1} - y_j}{x_{j+1} - x_j} \right| \approx \int \log |f'(x)| \, p(x) dx$$

 $\square$  infer  $X \rightarrow Y$  whenever

$$\hat{C}_{X \to Y} < \hat{C}_{Y \to X}$$

Janzing D, Mooij J, Zhang K, et al. Information-geometric approach to inferring causal directions. Artificial Intelligence, 2012.
### **Causal Function based method: practical issue - categorical variables**

DMIR

☑ Additive noise model for categorical variables

$$Y = g(X) + E, X \perp E$$



#### ☑ Problem of existing method: how to fit the discrete data? e.g. male, female, wood, iron, steel...

Cai R, Qiao J, Zhang K, et al. Causal discovery from discrete data using hidden compact representation. NeurIPS, 2018.

# **Causal Function based method: HCR**



☑ Categorical data: A Hidden Compact Representation Model



Cai R, Qiao J, Zhang K, et al. Causal discovery from discrete data using hidden compact representation. NeurIPS, 2018.

# **HCR: algorithm**



- ☑ Step 1: Estimate the model  $M: X \to Y' \to Y$ ,  $\widehat{M}: Y \to X' \to X$  by maximizing BIC  $L^*(M; D), L^*(\widehat{M}; D)$  respectively.
- $\square$  Step 2: If  $L^*(M; D) > L^*(\widehat{M}; D)$ , infer " $X \to Y$ "

If  $L^*(M; D) < L^*(\widehat{M}; D)$ , infer " $X \to Y$ "

If  $L^*(M; D) = L^*(\widehat{M}; D)$ , infer "non – identifiable"



# **HCR: Identifiability**



**Theorem 2**. Assume that in the causal direction there exists the transformation Y' = f(X) such that P(Y | X) = P(Y | Y'), where |Y'| < |X|, and assumption A1 holds. Then to produce the same distribution P(X, Y), the reverse direction must involve more effective number of parameters in the model than the causal direction.

X	P(X)	Y' Y	Food Poisoning		Stomac Flu	ch	Normal			
Poisonous Mushroom Mushroom	0.1	Poisonous Mushroom	0.85		0.10		0.05		Poisonous	
Rice	0.6	Mushroom	0.03	0.03			0.90	Not Poisonous		
Poisonous Fish	0.1	Rice	0.03		0.07		0.90		Not Poisonous	
		Poisonous Fish	0.85		0.10		0.05	Γ	Poisonous	
Y Food Poisoning	P(Y)	X' X	Poisonous Mushroom	Mus	hroom	Pois	sonous Fish			
Stomach Flu	0.079	Poisonous Mushroom	0.308	0.12	27	0.0	08	I	Poisonous Mushroom	
Normal	0.645	Mushroom	0.616	0.25	53	0.0	16		Mushroom	
		Rice	0.065	0.53	32	0.8	37		Rice	
		Poisonous Fish	0.011	0.08	39	0.1	40		Poisonous Fish	

Cai R, Qiao J, Zhang K, et al. Causal discovery from discrete data using hidden compact representation. NeurIPS, 2018.

# **HCR: results**



#### $\square$ The causal mechanism behind the abalone data set: Adult/Infant $\rightarrow$ size



#### Result on Abalone data set.



- ☑ 基于约束的因果发现方法
- ☑ 基于函数的因果发现方法
- ☑ 混合型因果发现方法

# How to handle the high dimensional data?



☑ Motivation: complimentary of Constraints based methods and Causal function based methods

	Causal Function based	Constraints based
High Dimensionality	NO (Pairwise)	Yes
Discovery Ability	Yes	NO (Markov Equivalence)

## **SELF: model**



☑ Embedding the functional causal assumption into the likelihood framework



$$\Pr(X_i = o_{j,i} \mid X_{P_i} = o_{j,P_i})$$

$$\underline{X_i = F_i(X_{P_i}) + E_i} \Pr(E_i = o_{j,i} - F_i(o_{j,P_i})|X_{P_i})$$

$$\underbrace{X_{P_i} \amalg E_i}_{\text{minimized}} \Pr\left(E_i = o_{j,i} - F_i(o_{j,P_i})\right)$$

Cai R, Qiao J, Zhang Z, et al. Self: structural equational likelihood framework for causal discovery. AAAI. 2018.

# **SELF: framework**



☑ Embedding the functional causal assumption into the likelihood framework



# **SELF: algorithm**



- $\square$  Step 1: For each nodes, fit  $F_i$  with a nonlinear or linear regression (e.g. XGBoost, OLS)
- ☑ Step 2: Estimate the noise distribution with kernel density function
- ☑ Step 3: Maximize



# **SELF: results**



#### $\ensuremath{\boxtimes}$ SOTA results, out of box

	F	1	Rec	call	Precision										
Data	SELF	ANM	SELF	ANM	SELF	ANM									
Child Alarm Win95pts Pathfinder	0.71 0.79 0.77	0.26 0.53 0.47 0.15	0.60 0.74 0.71	0.40 0.59 0.49	0.88 0.85 0.86	0.19 0.48 0.45 1.00									

Table 3: Results on real world structure with nonlinear data.

Table 2: Results on real world structure with linear data.

		I	71			Re	call		Precision			
Dataset	SELF	LiNGAM	DLiNGAM	HCBN	SELF	LiNGAM	DLiNGAM	HCBN	SELF	LiNGAM	DLiNGAM	HCBN
Child	0.98	0.98	0.95	0.58	1.00	1.00	1.00	0.65	0.96	0.95	0.92	0.52
Alarm	0.98	0.43	0.94	0.52	0.99	0.76	1.00	0.64	0.96	0.31	0.88	0.44
Win95pts	0.95	0.56	0.88	0.80	0.97	0.88	1.00	0.91	0.93	0.42	0.79	0.71
Pathfinder	0.91	0.86	0.85	0.73	0.95	0.96	0.96	0.83	0.87	0.77	0.76	0.64

**Cai R**, Qiao J, Zhang Z, et al. *Self: structural equational likelihood framework for causal discovery*. AAAI. 2018. Codes: <u>https://github.com/DMIRLAB-Group/CDMIR</u>



- ☑ 隐变量问题
- ☑ 非独立同分布问题

# **Causal discovery among latent variables**



- ☑ Discovering causal relations among latent variables is important in many domains
  - In psychotherapy: what is the causal relations among role conflict, depersonalization, personal accomplishment?





*Q***<sub>1</sub>**: What is the level of stress you are experiencing?

**Q**<sub>2</sub>: How often do you feel depressed?



#### ☑ Linear latent variable model

- Measurement model
- Structural model



Observed data

Causal structure of latent variables

# **Tetrad Condition**



**Tetrad Condition**: for  $X_1$ ,  $X_2$  and any one of  $X_3$ ,  $X_4$ , three quadratic constraints (tetrad constraints) on the covariance matrix are implied: e.g., for  $X_4$ ,

 $\rho_{12}\rho_{34} = \rho_{14}\rho_{23} = \rho_{13}\rho_{24},$ 

where  $\rho_{12}$  is the correlation between  $X_1, X_2$ , etc.

(Note that any two of the three vanishing tetrad differences above entails the third.)



Silva R, Scheines R, Glymour C, Spirtes P, Chickering DM. Learning the structure of linear latent variable models, JMLR 2006.



 $\square$  Cannot *identify* the causal direction:  $L_1 \rightarrow L_2$  or  $L_2 \rightarrow L_1$ ?



Cannot detect the latent variables when we only have 3 measured variables



How to solve this problem?

Intuition



☑ What kind of information is helpful to distinguish them?



Non-Gaussian information may help



☑ Consider the following two causal structures:



Asymmetry



$$Cov(X_i, X_k) = abd$$

$$Cov(X_j, X_k) = (b^2 + 1)cd$$

$$\frac{Cov(X_i, X_k)}{Cov(X_j, X_k)} = \frac{ab}{(b^2 + 1)c}$$

$$(X_i - \frac{Cov(X_i, X_k)}{Cov(X_j, X_k)} X_j) \searrow X_k$$



# **Triad condition**



**Triad Constraints (our proposed)** In a linear latent model, suppose  $\{X_i, X_j\}$  and  $X_k$  are distinct and correlated variables and that all noise variables are non-Gaussian. Define the pseudo-residual of  $\{X_i, X_j\}$  relative to  $X_k$ , which is called a reference variable, as

$$E_{(i,j|k)} := X_i - \frac{Cov(X_i, X_k)}{Cov(X_j, X_k)} \cdot X_j.$$

☑ We say that  $\{X_i, X_j\}$  and  $X_k$  satisfy Triad constraint if and only if  $E_{(i,j|k)} \perp X_k$ , i.e.,  $\{X_i, X_j\}$  and  $X_k$  violate the Triad constraint if and only if  $E_{(i,j|k)} \perp X_k$ .



# **Triad Condition: Limitation**



☑ Cannot applied into the multiple shared latent variables case:



☑ Can we extend the surrogate strategy to multiple shared latent variables?

# **Generalized Independent Noise (GIN) Condition**





 $E_{\mathbf{Y}}$  is independent from  $L_1$  and  $L_2$ .





#### Can we still use some measured variables as surrogate variables?

Similarly, use Measured Variables  $\mathbf{Z} = (X_4, X_5)^T$  as Surrogate Variables, we have

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = \begin{bmatrix} a_1 & b_1 \\ a_2 & b_2 \\ a_3 & b_3 \end{bmatrix} \cdot \frac{1}{a_5 b_4 - a_4 b_5} \begin{bmatrix} b_5 & b_4 \\ a_5 & a_4 \end{bmatrix} \begin{bmatrix} X_4 \\ X_5 \end{bmatrix} - \begin{bmatrix} a_1 & b_1 \\ a_2 & b_2 \\ a_3 & b_3 \end{bmatrix} \cdot \frac{1}{a_5 b_4 - a_4 b_5} \begin{bmatrix} b_5 & b_4 \\ a_5 & a_4 \end{bmatrix} \begin{bmatrix} \varepsilon_4 \\ \varepsilon_5 \end{bmatrix} + \begin{bmatrix} \varepsilon_{X_1} \\ \varepsilon_{X_2} \\ \varepsilon_{X_3} \end{bmatrix}$$

 $\exists$  nonzero vector  $\boldsymbol{\omega}$  s.t.  $\boldsymbol{\omega}^{\mathrm{T}} \mathbf{Cov}(\boldsymbol{Y}, \boldsymbol{Z}) = \mathbf{0}$ 

 $\boldsymbol{\omega} = [a_2b_3 - b_2a_3, b_1a_3 - a_1b_3, a_1b_2 - b_1a_2]^T$ 

 $\Rightarrow \boldsymbol{\omega}^{\mathrm{T}} \boldsymbol{Y}$  is independent from *L*.

Xie, F., Cai, R., Huang, B., Glymour, C., Hao, Z., & Zhang, K. Generalized Independent Noise Condition for Estimating Latent Variable Causal Graphs. NeurIPS 2020.

# **GIN: Generalized Independent Noise Condition**









 $\exists$  nonzero vector  $\boldsymbol{\omega}$  s.t. $\boldsymbol{\omega}^{\mathrm{T}} \cdot \mathbf{Cov}(\boldsymbol{Y}, \boldsymbol{Z}) = \mathbf{0} \Rightarrow \boldsymbol{\omega}^{\mathrm{T}} \boldsymbol{Y}$  is independent from *L*.  $\boldsymbol{\omega} = [a_2b_3 - b_2a_3, b_1a_3 - a_1b_3, a_1b_2 - b_1a_2]^T$ 

 $a_{3}a_{4} + b_{3}b_{4} \quad a_{3}a_{5} + b_{3}b_{4}$ 

#### GIN: view measured and latent variables in a unified way

Xie, F., Cai, R., Huang, B., Glymour, C., Hao, Z., & Zhang, K. Generalized Independent Noise Condition for Estimating Latent Variable Causal Graphs. NeurIPS 2020.

# **GIN: Generalized Independent Noise Condition**



Generalized Independent Noise(GIN) condition: (Z, Y) follows the GIN condition iff there

exists non-zeros  $\omega$  such that  $\omega^T \mathbb{E}[\mathbf{Y}\mathbf{Z}^T] = 0$  and  $\omega^T \mathbf{Y}$  is independent from  $\mathbf{Z}$ .

- $\blacksquare$  Triad condition can be seen as a special case of the GIN condition.
  - For example,  $(\{X_k\}, \{X_i, X_j\})$  satisfy Triad condition iff  $(\{X_k\}, \{X_i, X_j\})$  satisfy GIN condition



Xie, F., Cai, R., Huang, B., Glymour, C., Hao, Z., & Zhang, K. Generalized Independent Noise Condition for Estimating Latent Variable Causal Graphs. NeurIPS 2020.



### Linear Non-Gaussian Latent Variable Model (LiNGLaM)

A1. [Measurement Assumption] There is no observed variable in X being an ancestor of any latent variables in L.
 A2. [Non-Gaussianity Assumption] The noise terms are non-Gaussian.

A3. [Double-Pure Child Variable Assumption] Each

latent variable set  $\mathbf{L}'$ , in which every latent variable directly causes the same set of observed variables, has at least

 $2\text{Dim}(\mathbf{L}')$  pure measurement variables as children.

☑ A4. [Purity Assumption] There is no direct edge between observed variables.



A simple structure that satisfies LiNGLaM

# **GIN: Application in LiNGLaM**



- $\square$  We proposed a *two-steps* algorithm.
  - Step 1: find *causal clusters* (variables sharing the same latent variables as parents);
  - Step 2: determine *causal order* of the latent variables;
  - Estimate the coefficients if needed

# **GIN: Application in LiNGLaM**



#### $\square$ We proposed a *two-steps* algorithm.

- Step 1: find *causal clusters* (variables sharing the same latent variables as parents);
- Step 2: determine *causal order* of the latent variables;
- Estimate the coefficients if needed



Ground-truth graph



E.g., Test |latent|=1, we have

 $(\{X_1, ..., X_4, X_7, X_8\}, \{X_5, X_6\})$  satisfies GIN. Thus,  $\{X_5, X_6\}$  is a cluster.

Similarly,

 $(\{X_1, ..., X_4, X_5, X_6\}, \{X_7, X_8\})$  satisfies GIN. Thus,  $\{X_7, X_8\}$  is a cluster.

### The variables in a cluster share the same GIN conditions.

# **GIN: Application in LiNGLaM**



#### $\square$ We proposed a *two-steps* algorithm.

- Step 1: find *causal clusters* (variables sharing the same latent variables as parents);
- Step 2: determine *causal order* of the latent variables;
- Estimate the coefficients if needed



#### 64

# **GIN: Application in LiNGLaM**

GIN

Algorithm

- ☑ We simulate data following the LiNGLaM, including 4 cases, with different DAG structures for and measurement variables and latent variables.
- ☑ Goal: find clusters (determine the location of latent variables)?
  - Latent oimission: measure omitted latent variables

LSTC

• Latent commission: measure falsely detected latent variables

Latent omission

• Mismeasurements: measure the misclassification of observed variables

FOFC

BPC

<i>c</i>													
	500	0.00(0)	0.00(0)	1.00(10)	0.50(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)
Case 1	1000	0.00(0)	0.00(0)	1.00(10)	0.50(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)
	2000	0.00(0)	0.00(0)	1.00(10)	0.50(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)
	500	0.10(2)	0.20(4)	0.9(10)	0.50(10)	0.00(0)	0.05(1)	0.00(0)	0.00(0)	0.12(2)	0.12(4)	0.00(0)	0.20(10)
Case 2	1000	0.05(1)	0.15(3)	1.00(10)	0.50(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.04(1)	0.12(3)	0.00(0)	0.20(10)
	2000	0.00(0)	0.00(0)	1.00(10)	0.50(10)	0.00(0)	0.02(2)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.20(10)
	500	0.20(3)	0.20(3)	0.13(9)	0.10(1)	0.00(0)	0.03(3)	0.00(0)	0.00(0)	0.19(3)	0.17(3)	0.00(0)	0.00(0)
Case 3	1000	0.06(2)	0.13(2)	0.16(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.06(2)	0.00(0)	0.00(0)	0.00(0)
	2000	0.00(0)	0.00(0)	0.50(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.00(0)
	500	0.13(4)	0.40(6)	0.90(10)	0.63(10)	0.00(0)	0.23(5)	0.00(0)	0.00(0)	0.04(2)	0.15(6)	0.02(2)	0.06(4)
Case 4	1000	0.10(3)	0.26(6)	0.93(10)	0.66(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.05(3)	0.11(2)	0.01(1)	0.02(2)
	2000	0.03(1)	0.32(6)	1.00(10)	0.70(10)	0.00(0)	0.00(0)	0.00(0)	0.00(0)	0.04(1)	0.11(3)	0.00(10)	0.00(0)
Not	Note: The number in perentheses indicates the number of ecourter ages that the surrent algorithm econor												

Latent commission

FOFC

LSTC

BPC

GIN

Note: The number in parentheses indicates the number of occurrences that the current algorithm *cannot* correctly solve the problem.

#### Our proposed algorithm is more efficient and can find all latent variables!

GIN



FOFC

BPC

Mismeasurements

LSTC



# **GIN: Results**



- ☑ We apply our algorithm to discover the underlying causal structure behind the Teacher's Burnout, and we are developing a tool with the help of psychologist.
  - Discovered clusters and causal order of the latent variables:

Causal Clusters	Observed variables
$\mathcal{S}_{1}\left(1 ight)$	$RC_1, RC_2, WO_1, WO_2,$
	$DM_1, DM_2$
$\mathcal{S}_{2}\left(1 ight)$	$CC_1, CC_2, CC_3, CC_4$
$\mathcal{S}_{3}\left(1 ight)$	$PS_1, PS_2$
$\mathcal{S}_{4}\left(1 ight)$	$ELC_1, ELC_2, ELC_3, ELC_4,$
	$ELC_5$
$\mathcal{S}_{5}(2)$	$SE_1, SE_2, SE_3, EE_1,$
	$EE_2, EE_3, DP_1, PA_3$
$\mathcal{S}_{6}(3)$	$DP_2, PA_1, PA_2$

 $L(\mathcal{S}_1) > L(\mathcal{S}_2) > L(\mathcal{S}_3) > L(\mathcal{S}_5) > L(\mathcal{S}_4) > L(\mathcal{S}_6).$ 

(from root to leaf)



Hypothesized model by experts [Byrne, 2010]

Xie, F., Cai, R., Huang, B., Glymour, C., Hao, Z., & Zhang, K. Generalized Independent Noise Condition for Estimating Latent Variable Causal Graphs. NeurIPS 2020.

# **FRITL: More general setting**



☑ How to learn causal structure with latent variables from observed data?



#### **☑** Challenges:

- how to efficiently decompose a large global graph into local small structures *without introducing new latent confounders*?
- how to recover local structures accurately in the presence of latent confounders?

Chen W, Zhang K, Cai R, et al. FRITL: A Hybrid Method for Causal Discovery in the Presence of Latent Confounders. Submitted to JMLR, arXiv preprint arXiv:2103.14238, 2021.

# **FRITL: Algorithm**



- ☑ Step 1: Construct a PAG: based on conditional independence test
- ☑ Step 2: Infer local causal structures: IN condition
- ☑ Step 3: Detect shared latent confounders
- ☑ Estimate remaining undetermined local causal structures if needed





### ☑ 隐变量问题

### ☑ 非独立同分布问题

- Multi-domain
- Multi closely related domain

# How about data from multiple environments?



If data are generated from multiple environments, how to discover the causal structures from data?
 Data is non-i.i.d.



#### ☑ Challenges:

- how to group the subjects that are implied the same causal structure?
- how to recover the causal structure?

# **CCSL: Causal Clustering Structure Learning**





Chen, W., Wu, Y., Cai, R., et al. CCSL: A Causal Structure Learning Method from Multiple Unknown Environments. arXiv, 2022.

# **CCSL: causal clustering**



☑ Challenge 1: how to group the subjects that are implied the same causal structure?



Chen, W., Wu, Y., Cai, R., et al. CCSL: A Causal Structure Learning Method from Multiple Unknown Environments. arXiv, 2022.

# **CCSL: causal structural learning**





Chen, W., Wu, Y., Cai, R., et al. CCSL: A Causal Structure Learning Method from Multiple Unknown Environments. arXiv, 2022.
### **MD-LiNA: Multi-Domain causal discovery**



☑ Common Space V.S. Multi-domain space



Zeng, Y., Shimizu, S., Cai, R., Xie, F., Yamamoto, M., & Hao, Z. (2020). Causal discovery with multi-domain LiNGAM for latent factors. IJCAI 2021.

## **MD-LiNA: Formalization**



☑ The formalization of our model is shown in these three equations. To integrate the multi-domain data, we use a simple coding representation method.



$$f^{(m)} = \boldsymbol{B}^{(m)} \boldsymbol{f}^{(m)} + \boldsymbol{\varepsilon}^{(m)}$$

## **MD-LiNA: Formalization**



☑ The formalization of our model is shown in these three equations. To integrate the multi-domain data, we use a simple coding representation method.



 $\boldsymbol{f}^{(m)} = \boldsymbol{B}^{(m)} \boldsymbol{f}^{(m)} + \boldsymbol{\varepsilon}^{(m)}$ 

```
\Box \hspace{-1.5cm} \begin{array}{c} \bar{f} = H\tilde{f} \end{array}
```

Zeng, Y., Shimizu, S., Cai, R., Xie, F., Yamamoto, M., & Hao, Z. (2020). Causal discovery with multi-domain LiNGAM for latent factors. IJCAI 2021.

## **MD-LiNA: Formalization**



☑ The formalization of our model is shown in these three equations. To integrate the multi-domain data, we use a simple coding representation method.



Zeng, Y., Shimizu, S., Cai, R., Xie, F., Yamamoto, M., & Hao, Z. (2020). Causal discovery with multi-domain LiNGAM for latent factors. IJCAI 2021.



☑ 根因故障定位



☑ The data are generated from multiple topologically related environments, how to discover the causal structures from data?



☑ It is crucial to consider the topological structure behind the data for learning Granger causality.

Cai R, Wu S et al. THP: Topological Hawkes Processes for Learning Granger Causality on Event Sequences. Submitted to TNNLS, arXiv:2105.10884, 2021.

### **THP: The topological-temporal Hawkes model**

✓ How to learn the Granger causality among event types using the event sequences that generated by nodes in a topological network?



#### topological-temporal Hawkes model

Cai R, Wu S et al. THP: Topological Hawkes Processes for Learning Granger Causality on Event Sequences. Submitted to TNNLS, arXiv:2105.10884, 2021.

DMIR

Look at the Hawkes process from a temporal convolution perspective:

$$\lambda_{v}(t) = \mu_{v} + \sum_{v' \in \mathbf{PA}_{v}} \int_{t' \in \mathbf{T}_{t-}} \phi_{v',v}(t-t') dC_{v'}(t') \quad \longleftrightarrow \quad \lambda_{v}(t) = \mu_{v} + \sum_{v' \in \mathbf{PA}_{v}} (\phi_{v',v} * dC_{v'})_{\mathbf{T}}(t)$$

☑ Look at the topological structure from a graph convolution perspective :

$$y = g_{\theta}(L)s = g_{\theta}(U\Gamma U^{T})s = Ug_{\theta}(\Gamma)U^{T}s$$

 $\square$  Extend the causal intensity function to the topological-temporal domain  $G_N \times T$ .

$$\lambda_{v}(n,t) = \mu_{v} + \sum_{v' \in \mathbf{PA}} (\psi_{v',v} * dC_{v'})_{\mathcal{G}_{N} \times \mathbf{T}}(n,t),$$

Cai R, Wu S et al. THP: Topological Hawkes Processes for Learning Granger Causality on Event Sequences. Submitted to TNNLS, arXiv:2105.10884, 2021.

### The 1St in PCIC 2021 Causal Discovery challenge





Datasets: <u>https://competition.huaweicloud.com/information/1000041487/introduction</u> Codes: <u>https://github.com/DMIRLAB-Group/CDMIR</u>

# CAUSAL-LEARN:因果学习开源算法平台

## causal-learn: 因果学习算法平台



☑ 基于Python实现了经典和部分最新的因果学习算法。其中包含了因果发现的经典算法与API,并且提供了模块化的代码。



☆ » Welcome to causal-learn's documentation!

**O** Edit on GitHub

#### Welcome to causal-learn's documentation!

**causal-learn** is a Python translation and extension of the Tetrad java code. It offers the implementations of up-to-date causal discovery methods as well as *simple* and *intuitive* APIs.

#### Note

This project is under active development. For source code, please kindly refer to our GitHub Repository.

- GitHub: <u>https://github.com/cmu-phil/causal-learn</u>
- ☑ 文档: <u>https://causal-learn.readthedocs.io/en/latest/</u>
- ☑ 简单使用案例: <u>https://github.com/cmu-phil/causal-learn/tree/main/tests</u>



☑ causal-learn支持:

- 基于约束的因果发现方法(Constrained-based causal discovery methods): PC、FCI、CD-NOD算法等;
- 基于评分的因果发现方法(Score-based causal discovery methods): 包含BIC、BDeu、generalized score等评分的GES算法;
- 基于函数因果模型的因果发现方法(Functional causal models-based causal discovery methods):
   LiNGAM及其拓展方法、ANM、PNL等;
- 隐因果表征学习方法(Hidden causal representation learning): GIN方法;
- ・ 格兰杰因果分析 (Granger causal analysis);
- 多个独立的基础模块,比如独立性测试,评分函数,图操作,评测指标;
- 更多最新的因果发现算法,如gradient-based methods等。



### ☑ <del>安装</del>:可以通过pip来实现

pip install causal-learn

1 ! pip install causal-learn
Collecting causal-learn
Downloading causal_learn-0.1.2.3-py3-none-any.wh1 (163 kB)
163 kB 26.6 MB/s
Requirement already satisfied: statsmodels in /usr/local/lib/python3.7/dist-packages (from causal-learn) (0.10.2)
Requirement already satisfied: graphviz in /usr/local/lib/python3.7/dist-packages (from causal-learn) (0.10.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.7/dist-packages (from causal-learn) (1.0.2)
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from causal-learn) (4.64.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from causal-learn) (1.21.6)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.7/dist-packages (from causal-learn) (3.2.2)
Requirement already satisfied: pandas in /usr/local/lib/python3.7/dist-packages (from causal-learn) (1.3.5)
Requirement already satisfied: pydot in /usr/local/lib/python3.7/dist-packages (from causal-learn) (1.3.0)
Requirement already satisfied: networkx in /usr/local/lib/python3.7/dist-packages (from causal-learn) (2.6.3)
Requirement already satisfied: scipy in /usr/local/lib/python3.7/dist-packages (from causal-learn) (1.4.1)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.7/dist-packages (from matplotlib->causal-learn) (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib->causal-learn) (1.4.2)
Requirement already satisfied: pyparsing!=2.0.4, !=2.1.2, !=2.1.6, >=2.0.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib->causal-learn) (3.0.8)
Requirement already satisfied: python-dateutil>=2.1 in /usr/local/lib/python3.7/dist-packages (from matplotlib->causal-learn) (2.8.2)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/dist-packages (from kiwisolver>=1.0.1->matplotlib->causal-learn) (4.1.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (from python-dateutil>=2.1->matplotlib->causal-learn) (1.15.0)
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-packages (from pandas->causal-learn) (2022.1)
Requirement already satisfied: threadpoolct1>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-learn->causal-learn) (3.1.0)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn->causal-learn) (1.1.0)
Requirement already satisfied: patsy>=0.4.0 in /usr/local/lib/python3.7/dist-packages (from statsmodels->causal-learn) (0.5.2)
Installing collected packages: causal-learn
Successfully installed causal-learn-0.1.2.3



### ☑ 使用: causal-learn为所有模型都提供了简单易用的接口, 用户可以通过一行代码在自己的数据上进行因果发现

#### from causallearn.search.ConstraintBased.PC import pc

cg = pc(data, alpha, indep\_test, stable, uc\_rule, uc\_priority, mvpc, correction\_name, background\_knowledge, verbose, show\_progress)

*# visualization using pydot* 

cg.draw\_pydot\_graph()

```
3 dep 4 | (0, 2) with p-value 0.000000
3 dep 4 | (1, 2) with p-value 0.000000
Graph Nodes:
X1;X2;X3;X4;X5
Graph Edges:
1. X1 --- X2
2. X1 --- X3
3. X2 → X4
4. X3 → X4
5. X4 → X5
Depth=2, working on node 4: 100%]
```





### ☑ 使用: causal-learn为所有模型都提供了简单易用的接口, 用户可以通过一行代码在自己的数据上进行因果发现

```
from causallearn.search.FCMBased.GIN.GIN import GIN
G, K = GIN(data, indep test, alpha)
```

```
# visualization using pydot
pyd = GraphUtils.to_pydot(G)
pyd.write_png('test.png')
```

[[1, 2, 0], [8, 6, 7], [5, 3, 4]]

Runing the Step 1: Finding the Causal Clusters Runing the Step 2: Learning the Causal Order of Latent Variables Graph Nodes: L1;X2;X3;X1;L2;X9;X7;X8;L3;X6;X4;X5 Graph Edges: 1. L1  $\rightarrow$  X2 2. L1  $\rightarrow$  X3 3. L1  $\rightarrow$  X1 4. L1  $\rightarrow$  L2 5. L1  $\rightarrow$  L3 6. L2  $\rightarrow$  X9 7. L2  $\rightarrow$  X7 8. L2  $\rightarrow$  X8 9. L2  $\rightarrow$  L3 10. L3  $\rightarrow$  X6 11. L3  $\rightarrow$  X4 12. L3  $\rightarrow$  X5



### causal-learn: 使用方法



#### ☑ 可视化结果与评测: 在算法运行结束后,用户可以查看生成的因果图,并通过多种评测指标来与基准图进行对比

from causallearn.utils.GraphUtils import GraphUtils

# visualization using pydot
pyd = GraphUtils.to\_pydot(G)
pyd.write\_png('test.png')





### Causal discovery = Causal thinking/assumptions + ML/Statistical tools