

Causality and Machine Learning (80-816/516)

Classes 18 (March 20, 2025)

Causal Representation Learning 1: IID Case

Instructor:

Kun Zhang (kunzl@cmu.edu)

Zoom link: <u>https://cmu.zoom.us/j/8214572323</u>)

Office Hours: W 3:00-4:00PM (on Zoom or in person); other times by appointment

Uncover Causality from Observational Data: How?



- Causal discovery (Spirtes et al., 1993)/ causal representation learning (Schölkopf et al., 2021): find such representations with identifiability guarantees
- Causal system has "irrelevant" modules (Spirtes et al., 1993; Pearl, 2000)



- conditional independence among variables;
- independent noise condition;
- minimal (and independent) changes...

Footprint of causality in data

• Three dimensions of the problem:

i.i.d. data?	Parametric constraints?	Latent confounders?
Yes	No	No
No	Yes	Yes

We Mainly Focused on the IID Case: Recent Advances in Causal Representation Learning

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?	
		No	(Different types of)	PC ECL oto
	INO	Yes	equivalence class	
Yes		No	Unique identifiability	- LiNGAM
	Yes	Yes	(under structural conditions)	- Rank-based, GIN
		No	(Extended) regression	
Non-I, but I.D.	No/Yes	Yes	Latent temporal causal processes identifiable!	
	No	No	More informative than MEC (CD-NOD)	
I., but non-I.D.	Yes		May have unique identifiability	
	No	Vee	Changing subspace identifiable	
	Yes	res	Variables in changing relations identifiable	

CRL in IID Case: How to achieve it?

- Linear-Gaussian case
 - Tetrad conditions
 - Rank deficiency-based method
- Linear, non-Gaussian case
 - Theoretical results
 - GIN-based method
- Nonlinear case
 - Sparsity
- Summary: Why is it possible?

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?	
	A La	No	(Different types of)	
V	INO	Yes	equivalence class	
Yes	Vac	No	Unique identifiability	
	res	Yes	(under structural conditions)	

A Problem in Psychology: Finding Underlying Harametric Latent Mental Conditions?

i.i.d. data?	Parametric constraints?	Latent confounders?
Yes	No	No
No	Yes	Yes

• 50 questions for big 5 personality test

race	age	engnat	gender	hand	source	country	E1	E2	E3	E4	E5	E6	E7	E 8	E9	E10	N1	N2	N3	N4	N5	N6	N7	N8	N9	N10	A1	A2	A 3	A 4	A5
3	53	1	1	1	1	US	4	2	5	2	5	1	4	3	5	1	1	5	2	5	1	1	1	1	1	1	1	5	1	5	2
13	46	1	2	1	1	US	2	2	3	3	3	3	1	5	1	5	2	3	4	2	3	4	3	2	2	4	1	3	3	4	۷
1	14	2	2	1	1	PK	5	1	1	4	5	1	1	5	5	1	5	1	5	5	5	5	5	5	5	5	5	1	5	5	1
3	19	2	2	1	1	RO	2	5	2	4	3	4	3	4	4	5	5	4	4	2	4	5	5	5	4	5	2	5	4	4	3
11	25	2	2	1	2	US	3	1	3	3	3	1	3	1	3	5	3	3	3	4	3	3	3	3	3	4	5	5	3	5	1
13	31	1	2	1	2	US	1	5	2	4	1	3	2	4	1	5	1	5	4	5	1	4	4	1	5	2	2	2	3	4	3
5	20	1	2	1	5	US	5	1	5	1	5	1	5	4	4	1	2	4	2	4	2	2	3	2	2	2	5	5	1	5	1
4	23	2	1	1	2	IN	4	3	5	3	5	1	4	3	4	3	1	4	4	4	1	1	1	1	1	1	2	5	1	4	3
5	39	1	2	3	4	US	3	1	5	1	5	1	5	2	5	3	2	4	5	3	3	5	5	4	3	3	1	5	1	5	1
3	18	1	2	1	5	US	1	4	2	5	2	4	1	4	1	5	5	2	5	2	3	4	3	2	3	4	2	3	1	4	2
3	17	2	2	1	1	П	1	5	2	5	1	4	1	4	1	5	5	3	5	3	2	5	3	3	4	3	2	4	2	4	1
13	15	2	1	1	1	IN	3	3	5	3	3	3	2	4	3	3	1	5	3	3	2	3	2	3	2	4	4	4	2	2	5
13	22	1	2	1	2	US	3	3	4	2	4	2	2	3	4	3	3	3	3	3	2	2	4	4	2	3	1	4	1	5	1
3	21	1	2	1	5	US	1	3	2	5	1	1	1	5	1	5	5	3	5	2	5	5	3	2	5	3	1	1	1	4	2
3	28	2	2	1	2	US	3	3	3	4	3	2	2	4	3	5	2	4	4	4	4	4	2	2	3	2	1	4	2	4	2
3	21	1	1	1	5	US	2	3	2	3	3	1	1	3	4	4	2	4	2	4	1	2	2	2	2	2	4	2	4	2	5
13	19	1	2	1	2	FR	1	3	2	4	2	4	1	4	3	4	4	2	3	2	1	3	1	2	2	3	4	2	3	1	4
3	21	1	2	1	5	US	4	1	5	2	5	1	5 5	3	5	1	5	2	5	2	3	3	3	3	4	2	1	5	2	5	2

Learning Hidden Variables & Their Relations

i.i.d. data?	Parametric constraints?	Latent confounders?
Yes	No	No
No	Yes	Yes

• <u>Measured</u> variables (e.g., answer scores in psychometric questionnaires) were <u>generated by causally related latent variables</u>



• Find latent variables L_i and their causal relations from measured variables X_i ? 6

CRL in IID Case: How to achieve it?

- Linear-Gaussian case
 - Tetrad conditions
 - Rank deficiency-based method
- Linear, non-Gaussian case
 - Theoretical results
 - GIN-based method
- Nonlinear case
 - Sparsity
- Summary: Why is it possible?

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?	
	A La	No	(Different types of)	
Var	INO	Yes	equivalence class	
res	Vac	No	Unique identifiability	
	res	Yes	(under structural conditions)	



Tetrad Constraints $X_1 = \lambda_1 L + \varepsilon_1 (L)$ $X_2 = \lambda_2 L + \varepsilon_2$ $X_3 = \lambda_3 L + \varepsilon_3$ $X_1 X_2 X_3 X_4 X_1 Y_1 Y_2 Y_3 X_3 X_2 X_1 Y_1 Y_2$ $\sigma_{X_1X_2}\sigma_{X_3X_4} = (\lambda_1\lambda_2\sigma_L^2)(\lambda_3\lambda_4\sigma_L^2) = (\lambda_1\lambda_3\sigma_L^2)(\lambda_2\lambda_4\sigma_L^2) = \sigma_{X_1X_3}\sigma_{X_2X_4}$ $= (\lambda_1 \lambda_2 \sigma_I^2) (\lambda_3 \lambda_4 \sigma_I^2) = (\lambda_1 \lambda_4 \sigma_I^2) (\lambda_2 \lambda_3 \sigma_I^2) = \sigma_{X_1 X_4} \sigma_{X_2 X_3}$ (Intuition! Same for correlation coefficients)

Tetrad Constraints $X_1 = \lambda_1 L + \varepsilon_1 (L)$ $X_2 = \lambda_2 L + \varepsilon_2$ $X_3 = \lambda_3 L + \varepsilon_3$

Charles Spearman (1904)



Statistical Constraints \rightarrow

Measurement Model Structure



- Identify the structure from the constraint?
 - Assumptions (LMC, faithfulness, linearity)
 - In this case we can recover the structure if X_i are correlated

An illustration of Tetrad Conditions

Dep₁



 ρ_{ij} denotes the correlation coefficient between x_i and x_j

Applications of Tetrad Conditions

- One-factor measurement model [Silva et al., 2006, Kummerfeld et al., 2016]
- Tree structure [Pearl, 1988, Choi et al., 2011]

One-factor measurement model





CRL in IID Case: How to achieve it?

- Linear-Gaussian case
 - Tetrad conditions
 - Rank deficiency-based method
- Linear, non-Gaussian case
 - Theoretical results
 - GIN-based method
- Nonlinear case
 - Sparsity
- Summary: Why is it possible?

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?	
	A.(No	(Different types of)	
V	INO	Yes	equivalence class	
Yes	Vac	No	Unique identifiability	
	fes	Yes	(under structural conditions)	

t-separation

• Trek-separation

Definition 10.2 (Treks (Sullivant et al., 2010)). In a DAG \mathcal{G} , a trek from node X to node Y is an ordered pair of directed paths (P_1, P_2) where P_1 has a sink X, P_2 has a sink Y, and both P_1 and P_2 have the same source (the source of a directed path is the starting node from which the path originates).

Definition 10.3 (t-separation (Sullivant et al., 2010)). Let A, B, C_A, and C_B be four subsets of V in graph \mathcal{G} (not necessarily disjoint). (C_A,C_B) t-separates A from B if for every trek (P_1 , P_2) from a vertex in A to a vertex in B, either P_1 contains a vertex in C_A or P_2 contains a vertex in C_B.

- From the draft of "Causal Representation Learning" (distributed in class)

Why t-separation?

• t-separation and rank constraints:

Theorem 10.4 (Rank and t-separation (Sullivant et al., 2010)). Given two sets of variables A and B from a linear model with graph G, we have:

 $rank(\Sigma_{\mathbf{A},\mathbf{B}}) \le \min\{|\mathbf{C}_{\mathbf{A}}| + |\mathbf{C}_{\mathbf{B}}| : (\mathbf{C}_{\mathbf{A}},\mathbf{C}_{\mathbf{B}}) \text{ t-separates } \mathbf{A} \text{ from } \mathbf{B} \text{ in } \mathcal{G}\}, (10.2)$

where $\Sigma_{A,B}$ is the cross-covariance over A and B, and equality generically holds.

- Under the so-called rank faithfulness, the inequality relation in Theorem 10.4 becomes equality.
- Check the rank of sub-matrices of the covariance matrix and infer the minimal size of $|C_A|+|C_B|$ for A and B to be t-separated, and infer relations

Connection to d-separation

Theorem 10.5 (t- and d-sep (Di, 2009)). For disjoint sets A, B and C, C d-separates A and B in graph \mathcal{G} , iff there is a partition $\mathbf{C} = \mathbf{C}_{\mathbf{A}} \cup \mathbf{C}_{\mathbf{B}}$ such that $(\mathbf{C}_{\mathbf{A}}, \mathbf{C}_{\mathbf{B}})$ t-separates $\mathbf{A} \cup \mathbf{C}$ from $\mathbf{B} \cup \mathbf{C}$.

• But it can be more informative than d-separation:

EXAMPLE 2.13 (Spiders). Consider the graph in Figure 2 which we call a *spider*.

Clearly, we have that $(\{c\}, \{c\})$ t-separates A from B, so that the submatrix $\Sigma_{A,B}$ has rank at most 2. Although this rank condition must be implied by CI rank constraints on Σ and the fact that Σ is positive definite, it does not appear to be easily derivable from these constraints.



- S. Sullivant K. Talaska, and J. Draisma, "Trek separation for Gaussian graphical models", Annals of Statistics, 2008

Linear, Gaussian Case: With **Rank Deficiency** Constraints



- Can we find L_6 ?
 - $\Sigma_{(X_{10},X_{11}), \mathbf{X} \setminus \{X_{10},X_{11}\}} = 1$
- Recovering the equivalence class
 - With rank deficiency of crosscovariance matrices
 - recursively and cleverly

- Huang, Low, Xie, Glymour, Zhang, "Latent Hierarchical Causal Structure Discovery with Rank Constraints," NeurIPS 2022

Linear, Gaussian Case: With **Rank Deficiency** Constraints



- Can we find L_6 ?
 - $\Sigma_{(X_{10},X_{11}), \mathbf{X} \setminus \{X_{10},X_{11}\}} = 1$
- Recovering the equivalence class
 - With rank deficiency of crosscovariance matrices

- Conditional independence is a special case - $\operatorname{rank}(\Sigma_{(X1, X2), (X2, X3)}) = 1 \Leftrightarrow X_1 \parallel X_3 \mid X_2$

- Unified causal discovery based on rank deficiency constraints

- Dong, Huang, Ng, Song, Zheng, Jin, Legaspi, Spirtes, Zhang, "A Versatile Causal Discovery Framework to Allow Causally-Related Hidden Variables," ICLR 2024

Example: Big 5 Questions Are Well Designed but...

Big 5: openness; conscientiousness; extraversion; agreeableness; neuroticism



 Dong, Huang, Ng, Song, Zheng, Jin, Legaspi, Spirtes, Zhang, "A Versatile Causal Discovery Framework to Allow Causally-Related Hidden Variab²9s," ICLR 2024

Example: Big 5 Questions Are Well Designed but...



CRL in IID Case: How to achieve it?

- Linear-Gaussian case
 - Tetrad conditions
 - Rank deficiency-based method
- Linear, non-Gaussian case
 - Theoretical results
 - GIN-based method
- Nonlinear case
 - Sparsity
- Summary: Why is it possible?

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?	
	A Le	No	(Different types of)	
Ver	INO	Yes	equivalence class	
res	Vac	No	Unique identifiability	
	Tes	Yes	(under structural conditions)	

Necessary & Sufficient Conditions on the Structure: Linear, non-Gaussian case

i.i.d. data?	Parametric constraints?	Latent confounders?
Yes	No	No
No	Yes	Yes

- Allow a large number of latent variables
- Minimality has to be assumed
- Estimation is generally difficult

Identifiable graphs with only 3 measured variables



- Adams, Hansen, Zhang, "Identification of Partially Observed Linear Causal Models: Graphical Conditions for the Non-Gaussian and Heterogeneous Cases," NeurIPS 2021

CRL in IID Case: How to achieve it?

- Linear-Gaussian case
 - Tetrad conditions
 - Rank deficiency-based method
- Linear, non-Gaussian case
 - Theoretical results
 - GIN-based method
- Nonlinear case
 - Sparsity
- Summary: Why is it possible?

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?		
	A La	No	(Different types of)		
Ve -	INO	Yes	equivalence class		
Yes	Vac	No	Unique identifiability		
	res	Yes	(under structural conditions)		



• Find direction between latent variables L_1 and L_2 ?

- Xie, Cai, Huang, Glymour, Hao, Zhang, "Generalized Independent Noise Condition for Estimating Linear Non-Gaussian Latent Variable Causal Graphs," NeurIPS 2020
- Cai, Xie, Glymour, Hao, Zhang, "Triad Constraints for Learning Causal Structure of Latent Variables," NeurIPS 2019



 Xie, et al., "Generalized Independent Noise Condition for Estimating Linear Non-Gaussian Latent Variable Causal Graphs," NeurIPS 2020



Let $\mathbb{Z} = \{X_1\}$ and $\mathbb{Y} = \{X_2, X_3\}$, GIN!

- GIN condition: (**Z**, **Y**) follows GIN $\Leftrightarrow w^{\intercal}\mathbf{Y} \parallel \mathbf{Z}$ for nonzero w
 - has graphical implications

• Generalized Independent Noise (GIN) Condition:

(Z, Y) follows the GIN condition $\iff \omega^{\top} Y \perp Z$, where $\omega^{\top} \text{Cov}(Y, Z) = 0$ and $\omega \neq 0$

• Graphical criterion

(Z, Y) follows the GIN condition iff there is an exogenous set S of PA(Y) that blocks all paths between Y and Z, where $0 \le |S| \le \min(|Z|, |Y|-1)$



X_i: observed variables L_i: latent variables

- Xie, Cai, Huang, Glymour, Hao, Zhang, "Generalized Independent Noise Condition for Estimating Linear Non-Gaussian Latent Variable Causal Graphs," NeurIPS 2020
- Cai, Xie, Glymour, Hao, Zhang, "Triad Constraints for Learning Causal Structure of Latent Variables," NeurIPS 2019

GIN for Estimating Linear, Non-Gaussian LV Model

A two-step algorithm to identify the latent variable graph
By testing for GIN conditions over the input X₁, …, X₈



Step 2: determine *causal structure* of the latent variables



GIN-Based Method: Application to Teacher's Burnout Data

- Contains 28 measured variables
- 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

 $L(S_1) > L(S_2) > L(S_3) > L(S_5) > L(S_4) > L(S_6).$ (from root to leaf)

• Consistent with the hypothesized model





- Xie, Cai, Huang, Glymour, Hao, Zhang, "Generalized Independent Noise Condition for Estimating Linear Non-Gaussian Latent Variable Causal Graphs," NeurIPS 2020
- Cai, Xie, Glymour, Hao, Zhang, "Triad Constraints for Learning Causal Structure of Latent Variables," NeurIPS 2019

CRL in IID Case: How to achieve it?

- Linear-Gaussian case
 - Tetrad conditions
 - Rank deficiency-based method
- Linear, non-Gaussian case
 - Theoretical results
 - GIN-based method
- Nonlinear case
 - Sparsity
- Summary: Why is it possible?

i.i.d. data?	Parametric constraints?	Latent confounders?	What can we get?
	No Yes	No	(Different types of) equivalence class
		Yes	
Yes		No	Unique identifiability
		Yes	conditions)

Nonlinear Cases Generally Non-Identifiable

i.i.d. data?	Parametric constraints?	Latent confounders?	
Yes	No	No	
No	Yes	Yes	

• Nonlinear ICA:

- Generative model: $\mathbf{X} = f(\mathbf{S})$, where **S** has independent components
- De-mixing procedure: **Y** = g(**X**), where **Y** components are as independent as possible
- Solutions always exist and are highly non-unique: Why?

Identifiability of Nonlinear ICA: Structural Sparsity

i.i.d. data?	Parametric constraints?	Latent confounders?	
Yes	No	No	
No	Yes	Yes	



(Structural Sparsity) For all $k \in \{1, ..., n\}$, there exists C_k such that $\bigcap_{i \in C_k} \operatorname{supp}(\mathbf{J}_{\mathbf{f}}(\mathbf{s})_{i,:}) = \{k\}.$



- **Graphically**, for every latent variable S_i , there exists a set of observed variable(s) such that the intersection of their/its parent(s) is S_i
- **Example:** for S_1 , there exists X_1 and X_4 such that the intersection of their parents is S_1

- Zheng, Ng, Zhang, On the Identifiability of Nonlinear ICA: Sparsity and Beyond, NeurIPS 2022

Further Generalization of Nonlinear ICA

- Undercompleteness
 - More observed variables than latent variables
- Partial sparsity
 - Sparsity is violated for some variables
- Partial source dependence
 - Source independence violated for some variables
- Flexible grouping structures
 - Dependence within each group, independence across groups



Applied to real-world datasets (EMNIST)



Figure 4: Percentage of random structures satisfying Structural Sparsity w.r.t. different degree of undercompleteness (i.e., m/n).



Figure 5: Percentage of sources satisfying *Structural Sparsity* w.r.t. different numbers of sources in the bijective setting (m/n = 1).

- Zheng and Zhang, Generalizing Nonlinear ICA beyond Structural Sparsity, NeurIPS 2023 (oral)

Identifiability of nonlinear ICA: real-world images

Line thickness

Angle

Upper width

Height



Identification results on EMNIST

Each row represents an identified source with its value varying

Summary: CRL in IID Case

- Traditional causal discovery is a special case
- Why possible?
 - Sparsity! Stronger or weaker...
 - Linear-Gaussian case: strong
 - Nonparametric case (weak parametric constraints): strong
 - Linear, non-Gaussian case (additional parametric constraints): weaker
 - How to leverage non-IDD features of the data

Where Are We?

i.i.d. data?	Parametric constraint?	Latent confounders?	What can we get?
Yes	No	No	(Different types of) equivalence class
		Yes	
	Yes	No	Unique identifiability
		Yes	conditions)
Non-I, but I.D.	No/Yes	No	
		Yes	
I., but non-I.D.	No	NLa	5
	Yes	INO	
	No	Yes	
	Yes		