

On Learning Necessary and Sufficient Causal Graphs

Hengrui Cai, Yixin Wang, Michael Jordan, Rui Song (2023+, arXiv)

Presented by Rohit Kanrar

June 14, 2023

Table of Contents

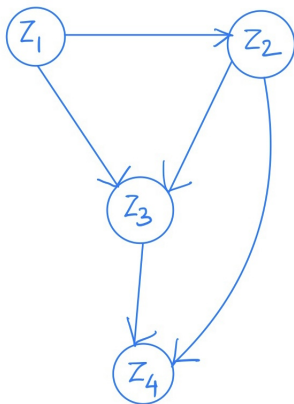
- 1 Preliminaries: Causal Discovery
- 2 Necessary & Sufficient Causal Graph (NSCG): Definition & Motivation
- 3 How to Quantify Necessity & Sufficiency in an NSCG?
- 4 NSCSL: A Causal Discovery Algorithm for Finding NSCG

Table of Contents

- 1 Preliminaries: Causal Discovery
- 2 Necessary & Sufficient Causal Graph (NSCG): Definition & Motivation
- 3 How to Quantify Necessity & Sufficiency in an NSCG?
- 4 NSCSL: A Causal Discovery Algorithm for Finding NSCG

Preliminaries: Causal Discovery

- Suppose we have p variables; $\mathbf{Z} = (Z_1, \dots, Z_p)$. Consider a Directed Acyclic Graph (DAG), $\mathcal{G} = (\mathbf{Z}, E)$, which has edge set E that contains only directed edges and no cyclic path.
- Causal Discovery (or Structure Learning) aims to identify $\hat{\mathcal{G}}$ based on the data, $\{\mathbf{Z}_i\}_{i=1}^n$ that closely resembles the true causal relationships between the variables, denoted by \mathcal{G} .



- Here we have 4 variables (or nodes).
- E contains 5 directed edges. There are no cyclic paths.
- If we include another edge, $Z_4 \rightarrow Z_1$, then it is no longer a DAG.
- Let $PA_{Z_i}(\mathcal{G})$ is the parent of Z_i in the DAG \mathcal{G} . For example, $PA_{Z_1}(\mathcal{G}) = \phi$ and $PA_{Z_3}(\mathcal{G}) = \{Z_1, Z_2\}$.

Figure: Example of a DAG

- **Bayesian Networks:** Consider a parametric density $f(\cdot)$ and a DAG $\mathcal{G} = (\mathbf{Z}, E)$. The density factorizes according to a DAG \mathcal{G} if there exists a set of parameter values $\Theta = \{\theta_1, \dots, \theta_p\}$ such that,

$$f(z_1, \dots, z_p) = \prod_{i=1}^p f_i(z_i | PA_{Z_i}(\mathcal{G}) = PA_{Z_i}(\mathcal{G}); \theta_i)$$

Such a pair (\mathcal{G}, Θ) is a Bayesian Network that defines the joint distribution.

- The distribution of \mathbf{Z} is *DAG-perfect* if there exists a DAG \mathcal{G} such that,
 - i (Markov) Every independence constraints encoded by \mathcal{G} holds in (\mathcal{G}, Θ) .
 - ii (Faithfulness) Every independence constraints encoded by (\mathcal{G}, Θ) holds in the DAG \mathcal{G} .

- **Linear Structural Equation Model (LSEM):** Suppose $\mathcal{G} = (\mathbf{Z}, E)$ is DAG-perfect. Under the LSEM assumption, the following holds,

$$\mathbf{Z} = B^T \mathbf{Z} + \epsilon$$

where $B = ((b_{i,j}))_{i,j=1}^p$ is called the *adjacency matrix*.

- For j -th variable, we have,

$$Z_j = b_{1,j}Z_1 + \dots + b_{p,j}Z_p + \epsilon_j$$

- $Z_i \rightarrow Z_j$ is present in E iff $b_{i,j} \neq 0$
- Without further assumption on ϵ , LSEM can be identified upto a markov equivalence class of multiple DAGs.
- For the rest of this discussion, we will assume ϵ is *multivariate Gaussian* (LSEM with Gaussian Error) to enforce the identifiability of DAG from the observational data.

- Let us consider a specific case where, we have p variables, $\{Z_1, \dots, Z_{p-1}, Y\}$, where Y is a response variable and is not expected to causally affect rest of the $p - 1$ variables.
- Identifying *full* causal graph based without this added information may lead to spurious causal effects.
- For LSEM with Gaussian Error, we may consider the following parameterization,

$$\begin{bmatrix} \mathbf{Z} \\ Y \end{bmatrix} = \begin{bmatrix} B_Z^T & 0 \\ \boldsymbol{\theta}^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{Z} \\ Y \end{bmatrix} + \begin{bmatrix} \epsilon_Z \\ \epsilon_Y \end{bmatrix}$$

where $\mathbf{Z} = (Z_1, \dots, Z_{p-1})^T$.

- However, it only leads to only *sufficient* causal graphs and include many spurious causal variables that does not affect the response.

Table of Contents

- 1 Preliminaries: Causal Discovery
- 2 Necessary & Sufficient Causal Graph (NSCG): Definition & Motivation
- 3 How to Quantify Necessity & Sufficiency in an NSCG?
- 4 NSCSL: A Causal Discovery Algorithm for Finding NSCG

Necessary and Sufficient Causal Graph (NSCG)

What is Sufficient Causal Graph? To present the idea, we first define some notations.

- $O = (\mathbf{Z}, Y)$ a collection of nodes containing a massive amount of features $\mathbf{Z} = [Z_1, \dots, Z_p]^T \in \mathcal{Z} \subset \mathbb{R}^p$ and the outcome of interest as $Y \in \mathcal{Y} \subset \mathbb{R}$.
- $Y(\mathbf{Z} = \mathbf{z})$ be the potential value of Y that would be observed after setting variable \mathbf{Z} as \mathbf{z} .
- Similarly, define the potential outcome $Y(Z_i = z_i)$ by setting individual variable Z_i as z_i , while keeping the rest unchanged.

- Let $\mathbf{X} \subset \mathcal{X} \in \mathbb{R}^d$ ($d \ll p$) be either subset of \mathbf{Z} or $\mathbf{X} = f(\mathbf{Z})$ that indeed captures the causal relationship between \mathbf{Z} and Y .
- Here $f : \mathbb{R}^p \rightarrow \mathbb{R}^d$ is an unknown function that selects a low-dimensional latent variable \mathbf{X} . Often, \mathbf{X} is treated as a low-dimensional representation of high-dimensional \mathbf{Z} .
- Denote the full graph, $\mathcal{G}_O = (O, e_O)$ and the sub-graph $\mathcal{G}_V = (V, e_V)$ with $O = (\mathbf{Z}, Y)$ and $V = (\mathbf{X}, Y)$ are the respective causal nodes and e_V, e_O are independent noise.
- Let $\mathbb{P}_{\mathcal{G}_V}$ and $\mathbb{P}_{\mathcal{G}_O}$ are the density functions under the respective causal graph. e.g., for LSEM with Gaussian error, $\mathbb{P}_{\mathcal{G}_V}$ and $\mathbb{P}_{\mathcal{G}_O}$ will be the Gaussian density function.

Sufficient Graph: \mathcal{G}_V is a sufficient causal graph to capture relationship among \mathbf{Z} and Y with $\mathbf{X} \subset \mathbf{Z}$ or $\mathbf{X} = f(\mathbf{Z})$ if

$$\begin{aligned} \mathbb{P}_{\mathcal{G}_V} \{Y | PA_Y(\mathcal{G}_V)\} &= \prod_{X_i \in PA_Y(\mathcal{G}_V)} \mathbb{P}_{\mathcal{G}_V} \{X_i | PA_{X_i}(\mathcal{G}_V)\} \\ &= \mathbb{P}_{\mathcal{G}_O} \{Y | PA_Y(\mathcal{G}_O)\} \prod_{Z_i \in PA_Y(\mathcal{G}_O)} \mathbb{P}_{\mathcal{G}_O} \{Z_i | PA_{Z_i}(\mathcal{G}_O)\} \end{aligned}$$

- Sufficient graphs may contain spurious relations that are not causally relevant to the outcome Y .
- Necessary graphs avoids such spurious relations.

Necessary Graph: \mathcal{G}_V is a necessary causal graph to capture the causal relationship among \mathbf{Z} and Y with $\mathbf{X} \subset \mathbf{Z}$ or $\mathbf{X} = f(\mathbf{Z})$ if for any subset $\mathbf{W} \subset \mathbf{X}$ or $\mathbf{W} = g(\mathbf{X})$, we have

$$\begin{aligned} & \mathbb{P}_{\mathcal{G}_V} \{Y | PA_Y(\mathcal{G}_V)\} \prod_{X_i \in PA_Y(\mathcal{G}_V)} \mathbb{P}_{\mathcal{G}_V} \{X_i | PA_{X_i}(\mathcal{G}_V)\} \\ & \neq \mathbb{P}_{\mathcal{G}_U} \{Y | PA_Y(\mathcal{G}_U)\} \prod_{W_i \in PA_Y(\mathcal{G}_U)} \mathbb{P}_{\mathcal{G}_U} \{W_i | PA_{W_i}(\mathcal{G}_U)\} \end{aligned}$$

where \mathcal{G}_U is the causal graph for causal nodes $U = (\mathbf{W}, Y)$.

- Now, we will discuss some motivating examples on why it is important to incorporate necessity in estimating causal graphs.

Y : Growth yield of yeast; Z : Gene expressions

$YMR105C \rightarrow YMR090W$ is a spurious relation that holds no causal impact on Y .

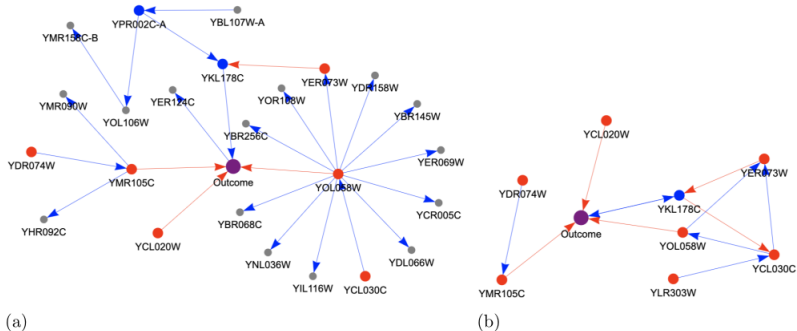


Figure: Causal graphs for candidate genes that affect the growth yield of yeast: (a) a sufficient graph; (b) a necessary and sufficient graph.

Y : Admission outcome of a student in a Graduate program.

Z : Other information of the student.

$G \rightarrow A$ is a spurious relation that holds no causal impact on Y .

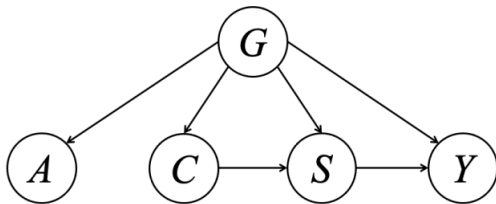


Figure: Causal relationships to understand the effect of gender in Graduate admissions. Node G defines the applicants' gender in the profile; node C is their pre-enrollment career objectives; node S is their choice of the department for study; node A is their appearance in the profile; node Y is the admission outcome.

Table of Contents

- 1 Preliminaries: Causal Discovery
- 2 Necessary & Sufficient Causal Graph (NSCG): Definition & Motivation
- 3 How to Quantify Necessity & Sufficiency in an NSCG?
- 4 NSCSL: A Causal Discovery Algorithm for Finding NSCG

How to find a Necessary and Sufficient Causal Graph (NSCG)?

- It is difficult to find an NSCG directly using the definitions, primarily due to the latent representation \mathbf{X} is unknown.
- A naive approach is to force search all different combinations of \mathbf{Z} to find a representation \mathbf{X} that satisfies the definitions; it is a nondeterministic polynomial (NP) hard problem.
- Instead, the authors suggest to evaluate the necessity and sufficiency of features on the prediction of outcome, using a property known as the probability of causation (POC).

Assumptions

Before presenting the definition of POC, we now present the assumptions needed for identification of causal quantities.

A1 Consistency: (or SUTVA)

$$\mathbf{Z} = \mathbf{z} \leftrightarrow Y(\mathbf{Z} = \mathbf{z}) = Y, \forall \mathbf{z} \in \mathcal{Z} \quad (1)$$

A2 Ignorability: (or NUC)

$$\begin{aligned} \{Y(\mathbf{Z} = \mathbf{z}), Y(\mathbf{Z} = \mathbf{z}')\} &\perp \mathbf{Z} \\ \{Y(Z_i = z_i), Y(Z_i = z'_i)\} &\perp Z_i | PA_{Z_i \cup Y}(\mathcal{G}_O) \end{aligned}$$

A3 Monotonicity:

$$\begin{aligned} \{Y(\mathbf{Z} \neq \mathbf{z}) = y\} \wedge \{Y(\mathbf{Z} = \mathbf{z}) \neq y\} &= False \\ \{Y(Z_i \neq z_i) = y\} \wedge \{Y(Z_i = z_i) \neq y\} &= False \end{aligned}$$

Probability of Causation (POC) in Tian and Pearl (2000)

POC for the feature \mathbf{Z} can be quantified using the probability of necessity and sufficiency.

Probability of Necessity and Sufficiency (PNS):

$$\begin{aligned} PNS &\stackrel{\text{def}}{=} \mathbb{P} \{ Y(\mathbf{Z} \neq \mathbf{z}) \neq y, Y(\mathbf{Z} = \mathbf{z}) = y \} \\ &\stackrel{\text{by (A1)}}{=} \mathbb{P} \{ \mathbf{Z} = \mathbf{z}, Y = y \} \quad PN + \mathbb{P} \{ \mathbf{Z} \neq \mathbf{z}, Y \neq y \} \quad PS \end{aligned}$$

where,

- the probability of necessity (PN) is defined as

$$PN \stackrel{\text{def}}{=} \mathbb{P} \{ Y(\mathbf{Z} \neq \mathbf{z}) \neq y | \mathbf{Z} = \mathbf{z}, Y = y \}$$

- the probability of sufficiency (PS) is defined as

$$PS \stackrel{\text{def}}{=} \mathbb{P} \{ Y(\mathbf{Z} = \mathbf{z}) = y | \mathbf{Z} \neq \mathbf{z}, Y \neq y \}$$

Intuition for PN and PS

Consider the feature $\mathbb{I}\{\mathbf{Z} = \mathbf{z}\}$ and the following notions,

$$\mathbf{Y} \begin{cases} = y & \text{is a "good" outcome} \\ \neq y & \text{is a "bad" outcome} \end{cases}$$

and

$$\mathbf{Z} \begin{cases} = \mathbf{z} & \text{is the feature being present} \\ \neq \mathbf{z} & \text{is the feature being absent} \end{cases}$$

- **Intuition of PN:** Given the good outcome observed (with the feature being present), PN measures the probability of a bad outcome after excluding the feature.

$$PN \stackrel{\text{def}}{=} \mathbb{P}\{Y(\mathbf{Z} \neq \mathbf{z}) \neq y | \mathbf{Z} = \mathbf{z}, Y = y\}$$

- **Intuition of PS:** Given the bad outcome observed (with the feature being absent) PS measures the probability of a good outcome after including the feature.

$$PS \stackrel{\text{def}}{=} \mathbb{P}\{Y(\mathbf{Z} = \mathbf{z}) = y | \mathbf{Z} \neq \mathbf{z}, Y \neq y\}$$

POC for Individual Feature Z_i

PNS can be further generalized to quantify the POC of an individual feature Z_i .

Conditional POC (C-POC):

$$C-POC_i \stackrel{\text{def}}{=} \mathbb{P} \{ Y(Z_i \neq z_i, \mathbf{Z}_{-i} = \mathbf{z}_{-i}) \neq \mathbf{y}, Y(Z_i \neq z_i, \mathbf{Z}_{-i} = \mathbf{z}_{-i}) = \mathbf{y} \}$$

Marginal POC (M-POC):

$$M-POC_i \stackrel{\text{def}}{=} \mathbb{P} \{ Y(Z_i \neq z_i) \neq \mathbf{y}, Y(Z_i \neq z_i) = \mathbf{y} \}$$

where, $\mathbf{Z}_{-i} \stackrel{\text{def}}{=} \mathbf{Z} \setminus Z_i$ be the complementary variable set of Z_i .

Theorem 4.4: Suppose (A1)-(A3) hold. Then the POC of an individual feature Z_i are all identifiable as,

$$C - POC_i = \mathbb{P}(Y = y | Z_i = z_i, \mathbf{Z}_{-i} \neq \mathbf{z}_{-i}) - \mathbb{P}(Y = y | Z_i \neq z_i, \mathbf{Z}_{-i} \neq \mathbf{z}_{-i})$$

$$M - POC_i = \mathbb{P}(Y = y | Z_i = z_i) - \mathbb{P}(Y = y | Z_i \neq z_i)$$

- Theorem 4.4 enables us to estimate POC from the observed data.
- But, estimating the conditional probabilities of Y based on high-dimensional features is challenging.
- This motivates the authors to consider expected mean outcome given different combinations of the confounders.

Corollary 4.5: Suppose (A1)-(A3) hold. Then we have

$$\begin{aligned}\int_{y \in \mathcal{Y}} y \, M - POC_i \, dy &= \mathbb{E} \{Y(Z_i = z_i)\} - \mathbb{E} \{Y(Z_i \neq z_i)\} \stackrel{\text{def}}{=} \delta_M(Z_i) \\ \int_{y \in \mathcal{Y}} y \, C - POC_i \, dy &= \mathbb{E} \{Y(Z_i = z_i, \mathbf{Z}_{-i} = \mathbf{z}_{-1})\} - \\ &\quad \mathbb{E} \{Y(Z_i \neq z_i, \mathbf{Z}_{-i} = \mathbf{z}_{-1})\} \\ &\stackrel{\text{def}}{=} \delta_C(Z_i)\end{aligned}$$

where $\delta_M(Z_i)$ and $\delta_C(Z_i)$ are defined as the marginal and conditional causal effects using the differences of expectations based on the corresponding POC.

Causal Effects for Non-binary Confounders

- The causal effects $\delta_M(\cdot)$ and $\delta_C(\cdot)$ are related to the POC only for binary confounders and positive outcome.
- To generalize the notion, authors extended the idea of total and direct effect for a variable of interest Z_i originally defined by Pearl (2009).

Natural Causal Effects:

$$TE_i = \partial \mathbb{E} \left\{ Y(Z_i = z'_i) \right\} / \partial z'_i$$

$$DE_i = \partial \mathbb{E} \left\{ Y(Z_i = z'_i, \mathbf{Z}_{-i} = \mathbf{z}_{-i}^{(z_i)}) \right\} / \partial z'_i$$

where $\mathbf{z}_{-i}^{(z_i)}$ is the value of \mathbf{Z}_{-i} if setting $do(Z_i = z_i)$.

Note: For binary confounders, the causal effects δ_M and δ_C match TE and DE , respectively.

Table of Contents

- 1 Preliminaries: Causal Discovery
- 2 Necessary & Sufficient Causal Graph (NSCG): Definition & Motivation
- 3 How to Quantify Necessity & Sufficiency in an NSCG?
- 4 NSCSL: A Causal Discovery Algorithm for Finding NSCG

Necessary and Sufficient Causal Structure Learning (NSCSL)

- A simple solution is to conduct a pre-screening process to find $\hat{\mathbf{X}} \subset \mathbf{Z}$, which achieve high C-POC/M-POC.
- Followed by employing a causal discovery method to obtain $\hat{\mathcal{G}} = \left(\left\{ \hat{\mathbf{X}}, Y \right\}, e \right)$ to approximate the NSCG, \mathcal{G}_V .
- Instead of such a two-step approach, the authors suggested a single-step procedure, NSCSL, which we will present now.

- Let $g : \mathbb{R}^p \rightarrow \mathbb{R}^d$ ($d \ll p$) be a selector function.
- We assume the nodes $\{g(\mathbf{Z}), Y\}$ have the following causal structure.

$$\begin{bmatrix} g(\mathbf{Z}) \\ Y \end{bmatrix} = \begin{bmatrix} B_Z^T & 0 \\ \boldsymbol{\theta}^T & 0 \end{bmatrix} \begin{bmatrix} g(\mathbf{Z}) \\ Y \end{bmatrix} + \begin{bmatrix} \epsilon_Z \\ \epsilon_Y \end{bmatrix}$$

where $\boldsymbol{\theta} \in \mathbb{R}^d$, $B_Z \in \mathbb{R}^{d \times d}$ and ϵ is a $d + 1$ dimensional random vector of joint independent error.

Note: This is the LSEM with Gaussian error assumption. Additionally, we impose the constraint of no descendent of Y .

How to estimate the POC for i -th selected feature $g_i(\mathbf{Z})$?

- Estimation of $\mathbb{P}(Y = y | g_i(\mathbf{Z}) = z_i)$ can be achieved by either parametric models (e.g., logistic for binary outcome) or non-parametric models (e.g., random forest/neural network)
- Data: $\{o^{(j)} = z^{(j)}, y^{(j)}\}_{1 \leq j \leq n}$.
- The marginal POC can be estimated as,

$$\widehat{M-POC} \left(g_i | \{o^{(j)}\} \right) = \prod_{j=1}^n \left| \widehat{\mathbb{P}} \left\{ Y = y^{(j)} | g_i(\mathbf{Z}) = g_i(\mathbf{z}^{(j)}) \right\} - \widehat{\mathbb{P}} \left\{ Y = y^{(j)} | g_i(\mathbf{Z}) \neq g_i(\mathbf{z}^{(j)}) \right\} \right|$$

- Similarly, the conditional POC can be estimated as,

$$\begin{aligned} & \widehat{C-POC} \left(g_i | \left\{ o^{(j)} \right\} \right) \\ &= \prod_{j=1}^n \left| \widehat{\mathbb{P}} \left\{ Y = y^{(j)} | g_i(\mathbf{Z}) = g_i(\mathbf{z}^{(j)}), g_{-i}(\mathbf{Z}) = g_{-i}(\mathbf{z}_{-i}^{(j)}) \right\} - \right. \\ & \quad \left. \widehat{\mathbb{P}} \left\{ Y = y^{(j)} | g_i(\mathbf{Z}) \neq g_i(\mathbf{z}^{(j)}), g_{-i}(\mathbf{Z}) = g_{-i}(\mathbf{z}_{-i}^{(j)}) \right\} \right| \end{aligned}$$

where $g_{-i}(\cdot) \stackrel{\text{def}}{=} g(\cdot) \setminus g_i(\cdot)$ is the complementary of $g_i(\cdot)$.

- Even for $d \ll p$, the estimation of M-POC/C-POC are difficult (Wang and Jordan (2022)).
- Authors proposed to use TE/DE as it has close form expression under LSEM.

Close Form Expressions of TE and DE under LSEM: Recall,

$$\begin{bmatrix} g(\mathbf{Z}) \\ Y \end{bmatrix} = \begin{bmatrix} B_Z^T & 0 \\ \boldsymbol{\theta}^T & 0 \end{bmatrix} \begin{bmatrix} g(\mathbf{Z}) \\ Y \end{bmatrix} + \begin{bmatrix} \epsilon_Z \\ \epsilon_Y \end{bmatrix}$$

- Then we have,

$$DE_i = \theta_i$$

as θ_i presents the weight of the direct edge $g_i(\mathbf{Z}) \rightarrow Y$.

- Total causal effect TE_i of the i -th selected feature $g_i(\mathbf{Z})$ can be expressed simplified using the path method, which we describe next.

- Let $\pi_i = \{g_i(\mathbf{Z}) \rightarrow \dots \rightarrow Y\}$ be the set of directed paths that starts with $g_i(\mathbf{Z})$ and ends with Y . Suppose there are m_i directed paths in π_i .
- The causal effect of $g_i(\mathbf{Z})$ on Y through the directed path $\pi_i^{(k)} = \{i, l_1, \dots, l_{e_k}, d+1\} \in \pi_i$ with length $e_k + 1$ is

$$PE \left\{ \pi_i^{(k)} \right\} = \prod b_{i, l_1} \dots b_{l_{e_k-1}, l_{e_k}} \theta_{l_{e_k}}$$

where $B_{\mathbf{Z}} = ((b_{i,j}))_{p \times p}$.

- Then, the total effect of $g_i(\mathbf{Z})$ on Y is defined as,

$$TE_i = \sum_{k=1}^{m_i} PE \left\{ \pi_i^{(k)} \right\}$$

NSCSL: Learning Algorithm

• Step 1: Causal Discovery

- Main idea is to estimate B with the acyclicity constraint.
- Authors proposed to use the following constraint proposed by Yu et al. (2019).

$$h_1(B) \stackrel{\text{def}}{=} \text{tr} [(I_{d+1} + tB \circ B)^{d+1}] - (d+1)$$

where $B \circ B$ is the element-wise square operator.

- Theorem 1 (Yu et al. (2019)) $\implies h_1(B) = 0$ iff B is acyclic.
- First part of the proposed loss,

$$L_1(B, g, \theta, \lambda_1 | \{o^{(j)}\}) = f(B, g, \theta | \{o^{(j)}\}) + \lambda_1 h_1(B)$$

where $f(\cdot)$ is any suitable loss function (such as least square for LSEM with Gaussian error).

• Step 2: Constraint based on POC/Natural Effects

- Based on POC:

$$L_2^P(B, g, \gamma | \{o^{(j)}\}) = - \sum_{i=1}^d \widehat{P}(g_i | \{o^{(j)}\}) + \gamma R(g)$$

where \widehat{P} is either $\widehat{C-POC}$ or $\widehat{M-POC}$.

- Based on Natural Effects:

$$L_2^{CE}(B, g, \gamma | \{o^{(j)}\}) = - \sum_{i=1}^d \widehat{CE}_i(B) + \gamma R(g)$$

where \widehat{CE}_i is either DE_i or TE_i .

- Here γ is a penalty term and $R(\cdot)$ is some norm to control the complexity of $g(\cdot)$.

- Finally, the proposed optimization problem is,

$$\min_{B, g} f \left(B, g, \theta \mid \left\{ o^{(j)} \right\} \right) + \lambda_1 h_1(B) + \gamma R(g)$$

- The estimated NSCG, \hat{g}_V can be obtained by using the estimated adjacency matrix, \hat{B} .
- Authors proposed to use a black-box stochastic optimizer, Adam for estimation purposes.

References

- Pearl, J. (2009). “Causal inference in statistics: An overview.”
- Tian, J. and Pearl, J. (2000). “Probabilities of causation: Bounds and identification.” *Annals of Mathematics and Artificial Intelligence*, **28(1-4)**, 287–313.
- Wang, Y. and Jordan, M.I. (2022). “Desiderata for representation learning: A causal perspective.”
- Yu, Y., Chen, J., Gao, T., and Yu, M. (2019). “Dag-gnn: Dag structure learning with graph neural networks.” In “International Conference on Machine Learning,” pages 7154–7163. PMLR.

Thank You