INDUCED COVARIANCE FOR CAUSAL DISCOVERY IN LINEAR SPARSE STRUCTURES

Anonymous authors

Paper under double-blind review

ABSTRACT

Causal models seek to unravel the cause-effect relationships among variables from observed data, as opposed to mere mappings among them, as traditional regression models do. This paper introduces Sparse Linear Causal Discovery (SLCD), a novel causal discovery algorithm designed for settings in which variables exhibit linearly sparse relationships. In such scenarios, the causal links represented by directed acyclic graphs (DAGs) can be encapsulated in a structural matrix. The proposed approach leverages the structural matrix's ability to reconstruct data and the statistical properties it imposes on the data to identify the correct structural matrix. This method does not rely on independence tests or graph fitting procedures, making it suitable for scenarios with limited training data. Simulation results demonstrate that the SLCD outperforms the well-known PC, GES, BIC exact search, and LINGAM-based methods in recovering linearly sparse causal structures by an average of 35% in precision and 41.5% in recall across all tested datasets.

1 Introduction

Causal learning is an approach used to extract and understand cause-and-effect relationships from data. This approach seeks to uncover the fundamental structures that determine how data are related Shaska & Mitra (2025). This structural understanding is at a deeper level than that observed in statistical learning, which is focused on learning various mappings among data Schölkopf & von Kügelgen (2022). Discovering causal relations plays a crucial role in the scientific method Camps-Valls et al. (2023). A comprehensive causal model of a phenomenon could describe the observed data and consistently make predictions. The advantage of this type of learning over statistical learning, which identifies mere associations between variables, lies in its generalization and robustness to distribution changes Schölkopf et al. (2021). Furthermore, causal relations may be transferable to other problems, which constitutes an additional benefit Schölkopf et al. (2021).

Several approaches have been proposed for causal discovery Pearl & Verma (1995); Spirtes et al. (2001); Shimizu et al. (2006; 2011); Meek (1997); Chickering (2002); Yuan & Malone (2013). These methods are generally classified into two categories Schölkopf & von Kügelgen (2022): constraint-based and score-based methods. In constraint-based methods, conditional independencies among variables are tested, and the inferred relations are represented using a DAG that best reflects them. Notable examples of such algorithms include inductive causation (IC) Pearl & Verma (1995), Spirtes-Glymour-Scheines (SGS) Spirtes et al. (2001), Peter-Clark (PC) Spirtes et al. (2001) and linear non-Gaussian acyclic model (LINGAM) based methods Shimizu et al. (2006) and Shimizu et al. (2011). IC and SGS algorithms examine the conditional independencies between each pair of variables conditioned on any subset of the remaining variables and use this information for forming the causal graph. This approach can be computationally intensive due to the large number of subsets. The PC algorithm mitigates this challenge by initiating the search from a complete graph and systematically removing edges, sequentially testing the conditional independencies of each pair and their neighbors. Although the PC algorithm reduces computational cost, it still depends on conditional independence tests, which are computationally demanding and require substantial data, particularly in high-dimensional settings, to produce reliable results. Unlike the previously mentioned methods that rely heavily on conditional independence tests, LiNGAM-based approaches assume non-Gaussianity and linearity in the data. By leveraging these assumptions, LiNGAM aims

to identify the causal graph. Although alleviating the challenge of conditional tests, non-Gaussianity is a limiting assumption.

Alternatively, score-based methods use a scoring function to evaluate graphical representations. Possible graphs are tested against the data, and the graph with the highest score is selected. Some of the prominent methods in this category are greedy equivalent search (GES) Meek (1997); Chickering (2002), and Bayesian information criterion (BIC) exact search Yuan & Malone (2013). The primary drawback of these methods Meek (1997); Chickering (2002); Yuan & Malone (2013) is the exponential growth of the possible graphs as the number of variables (nodes of the graph) increases, which results in higher computation demands.

Another line of work in causal discovery is *causal representation learning* Schölkopf et al. (2021); Varici et al. (2024). In this setting, data are assumed to be generated from high-level latent variables, which are mapped through transformations to the observed data. The objective is to identify both the relations among the latent variables and the transformations that connect them to the observed data. This is carried out in such a way that the resulting representation remains consistent with causal interventions.

However, these approaches do not adequately address scenarios with limited data, where conditional independence tests often fail to yield reliable results. Moreover, causal representation learning typically requires access to both the dataset and its intervened versions, which can be a significant limitation. Therefore, despite the emergence of several methods for causal discovery, there remains no algorithm well-suited for situations involving small datasets and the absence of feasible interventions.

The primary contribution of this paper is the development of a novel causal discovery algorithm designed for linear sparse structures. The key contributions are as follows:

- We propose *Sparse Linear Causal Discovery (SLCD)*, a new algorithm that recovers the structural matrix that encapsulates causal graph information by leveraging induced covariance, data reconstruction, rank, and diagonal structure, specifically for settings where variable relationships can be effectively modeled as sparse linear dependencies.
- We introduce the concept of *induced covariance*, a statistical property implied by causal structures, and provide its formal mathematical characterization.
- We extend the induced covariance framework to accommodate nonlinear causal structures.
- We provide a theoretical analysis of structural matrices that satisfy both induced covariance and reconstruction constraints, establishing results on local uniqueness and sensitivity.
- We demonstrate through experiments that the proposed method outperforms established causal discovery approaches, achieving on average a 35% improvement in precision and a 41.5% improvement in recall across all tested datasets.

The rest of the paper is organized as follows. Section 2 provides a brief overview of the causal learning framework. Section 3 outlines the problem statement and associated challenges, while Section 4 introduces the proposed approach (SLCD). Section 5 discusses the generalization of SLCD to settings with nonlinear causal relations. Section 6 presents the theoretical analysis of structural matrices that satisfy both induced covariance and reconstruction constraints. Section 7 presents and analyzes the simulation results, and Section 8 concludes the paper.

2 PRELIMINARIES

To state the problem, we first review the concept of a structural causal model (SCM), a popular method for causal relations modeling Schölkopf & von Kügelgen (2022). In this framework, a set of random variables $\{y_1, y_2, \ldots y_n\}$ is represented as the vertices of a directed acyclic graph (DAG), and the following relations hold:

$$y_i = f_i(\mathcal{P}_i, u_i) \quad \forall i \in \{1, 2, \dots, n\},\tag{1}$$

with f_i being a deterministic function, \mathcal{P}_i (parents) represents the set of variables that influence y_i , and u_i is an unexplained noise random variable Schölkopf et al. (2021). In the graphical SCM representation, a directed edge exists from each member of \mathcal{P}_i to y_i for $i \in \{1, 2, \dots, n\}$. The process of causal discovery involves identifying f_i and \mathcal{P}_i for all $n \in \{1, 2, \dots, n\}$.

3 PROBLEM STATEMENT

Let $\boldsymbol{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$ be a vector that contains all the random variables for which causal relationships are to be discovered. We define \mathcal{I} as the set of indices for independent variables and \mathcal{D} as the set of indices for dependent variables. According to the SCM model, each x_i for $i \in \mathcal{D}$ is a function of a subset of independent variables \mathcal{P}_i , which are considered the parents of x_i . We assume that functions $f_i, \forall i \in \{1, 2, \dots, n\}$ are linear and $|\mathcal{P}_i| \leq \tau$, where $|\cdot|$ represents the cardinality of a set and τ is a model parameter. Under these assumptions, we can represent x_i as follows:

$$x_i = \boldsymbol{d}_i^T \boldsymbol{x},\tag{2}$$

where $d_i \in \mathbb{R}^n$ is a vector with fewer than τ non-zero elements, corresponding to the independent variables upon which x_i depends. Given this notation, the relationship for all variables can be expressed as

$$x = Dx, (3)$$

where $D \in \mathbb{R}^{n \times n}$ is a matrix constructed from the vectors $d_i, \forall i \in \{1, 2, ..., n\}$ as its rows. The matrix D contains all pertinent information on the causal structure of this model, and an accurate estimation of D reveals the underlying causal relations.

In practice, there is often limited or no prior knowledge about the underlying structure of the data, and only the dataset itself is available. We use $X = [x_1, x_2, \dots x_m] \in \mathbb{R}^{n \times m}$ to represent the given dataset, where each $x_i \in \mathbb{R}^n$ is a sample. Applying (3), we have

$$X = DX. (4)$$

The primary objective is to determine the causal structure (D) from X. Estimation of D does not involve conditional independence tests, making it suitable when the number of data samples are limited, especially in high-dimensional data.

3.1 CHALLENGES

In the estimation of D, several challenges must be addressed. Based on the previous discussion, it can be inferred that D must satisfy the condition expressed in (3). However, as demonstrated in the following example, this condition alone is insufficient for uniquely determining the causal structure.

Example 1. Suppose data is created as follows

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} .$$
 (5)

This structure suggests that x_1 and x_2 are independent variables (as they are not linear combinations of any other variables), while x_3 is the sum of x_1 and x_2 . When only the data is available and the objective is to satisfy the condition given in (3), the solution may not be unique. For example, the identity matrix ($I \in \mathbb{R}^{3\times 3}$) and the following matrix also satisfies (3):

$$\begin{bmatrix} 0 & -1 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} . ag{6}$$

This is a common problem in causal discovery as multiple graphs can describe the same data Spirtes et al. (2001). This example also illustrates how causal discovery differs from a regression problem. While both solutions may be acceptable in the context of regression, only one solution reveals the underlying causal structure. In other words, what separates this approach from a linear algebra regression is that we are not looking for any solution of (3) but the one that describes the associations according to cause-and-effect relations.

4 INDUCED COVARIANCE-BASED CAUSAL DISCOVERY

To address the challenges discussed above, it is beneficial to explore certain properties D, which represents the causal relations.

The structure of D can provide valuable insights into the relations among variables. Any k^{th} row of the structural matrix D, whose all elements are all equal to zero except for the element in column k, which is equal to 1, corresponds to an independent random variable. Since all variables are linear combinations of these variables, the rank of D corresponds to the number of independent variables. This observation suggests that the structure of D can be used to infer the statistical relationships between the variables. This property is particularly useful for narrowing down the number of potential solutions for D.

To further constrain the possible solutions for D, the following theorem establishes a connection between the statistical properties of the data and the structural matrix. More specifically, it shows that selecting a specific value for variable D, uniquely determines the value of the covariance matrix of the data, indicating that D imposes a constraint on the covariance matrix resulting in *induced covariance*.

Theorem 1. Consider $D \in \mathbb{R}^{n \times n}$ to be a matrix that represents a linear causal structure governing the zero mean variables $\boldsymbol{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$. The covariance matrix of these variables is given by $D\boldsymbol{\sigma}D^T$, in which $\boldsymbol{\sigma} \in \mathbb{R}^{n \times n}$ is a diagonal matrix with diagonal elements being the variance of variables \boldsymbol{x} .

Proof. Let \mathbb{E} be the expected value operator. To prove this theorem, we derive an expression for $\mathbb{E}[x_i, x_j]$ where x_i and x_j are two components of the random vector \boldsymbol{x} . Based on the causal structure, the following holds

$$x_i = \boldsymbol{d}_i^T \boldsymbol{x},\tag{7}$$

$$x_j = \boldsymbol{d}_i^T \boldsymbol{x},\tag{8}$$

in which d_i^T and d_j^T are the rows i and j of the structural matrix D. Thus,

$$\mathbb{E}[x_i, x_j] = \mathbb{E}[d_i^T x d_j^T x]. \tag{9}$$

Since the only nonzero elements in $d_i^T x d_j^T x$ occurs when both d_j and d_i have non-zero elements in the same positions, then

$$\mathbb{E}[x_i, x_j] = \boldsymbol{d}_i^T \boldsymbol{\sigma} \boldsymbol{d}_j^T. \tag{10}$$

By applying the same procedure to all (i, j) pairs, the theorem is proven.

This theorem restricts the solutions of (3) by imposing that the correct solution must not only satisfy (3) but also fulfill the condition $\Sigma = D\sigma D^T$, where $D\sigma D^T$ is the induced covariance by D and $\Sigma \in \mathbb{R}^{n \times n}$ is the covariance matrix of data, which can be estimated directly from data.

By using the properties of D and its implications on the structure of data, we can formulate the following optimization problem for structure recovery:

$$\begin{aligned} \underset{\boldsymbol{D}}{\text{arg min}} \{ \text{rank}(\boldsymbol{D}) + \lambda \text{Tr}(\boldsymbol{D}) \} \\ \text{subject to} \quad \boldsymbol{X} &= \boldsymbol{D} \boldsymbol{X}, \\ \boldsymbol{\Sigma} &= \boldsymbol{D} \boldsymbol{\sigma} \boldsymbol{D}^T, \\ \|\boldsymbol{d}_i^T\|_0 &= \tau \quad \forall i \in \{1, 2, \dots, n\}. \end{aligned} \tag{11}$$

In this formulation, $D \in \mathbb{R}^{n \times n}$ represents the structural matrix, while $X = [x_1, x_2, \dots x_m] \in \mathbb{R}^{n \times m}$ is the dataset. The covariance matrix of the data is represented by $\Sigma \in \mathbb{R}^{n \times n}$, and $\sigma \in \mathbb{R}^{n \times n}$ is a diagonal matrix whose diagonal entries correspond to those of Σ . The term d_i^T represent row i of D. The operator $\|\cdot\|_0$ returns the number of non-zero elements in a vector. Additionally, τ controls the number of independent variables, and λ serves as a scaling parameter. The rank(D) term prevents the model from becoming overly complicated, and $\operatorname{tr}(D)$ discourages the solution from being close to the I, which implies all variables are independent.

Rank (the number of non-zero singular values) requires combinatorial calculation, which makes the problem untractable. To address this challenge, the idea proposed in Mohimani et al. (2009) is used, which approximate $\|.\|_0$ as:

$$||x||_0 \approx 1 - e^{-\frac{x^2}{\sigma^2}}. (12)$$

 By combining these ideas, the final problem formulation is

$$\underset{D}{\operatorname{arg\,min}} \{ \sum_{i=1}^{n} (1 - e^{-\frac{s_i^2}{\sigma^2}}) + \lambda \sum_{i=1}^{n} (1 - e^{-\frac{d_{(i,i)}^2}{\sigma^2}}) \}$$
subject to $\boldsymbol{X} = \boldsymbol{D}\boldsymbol{X}$,
$$\boldsymbol{\Sigma} = \boldsymbol{D}\boldsymbol{\sigma}\boldsymbol{D}^T,$$

$$\|\boldsymbol{d}_i^T\|_0 = \tau \quad \forall i \in \{1, 2, \dots, n\},$$
(13)

where $s_i, \forall i \in \{1, 2, ..., n\}$ are the singular values of \boldsymbol{D} .

To present the final algorithm for obtaining the solution of (13), it is necessary to consider $\|d_i^T\|_0 = \tau \quad \forall i \in \{1, \dots, n\}$, which also requires combinatorial calculations. To handle that, we propose solving the following optimization problem:

$$\underset{D}{\operatorname{arg\,min}} \{ \sum_{i=1}^{n} (1 - e^{-\frac{s_i^2}{\sigma^2}}) + \lambda \sum_{i=1}^{n} (1 - e^{-\frac{d_{(i,i)}^2}{\sigma^2}}) \}$$
subject to $\boldsymbol{X} = \boldsymbol{D}\boldsymbol{X}$,
$$\boldsymbol{\Sigma} = \boldsymbol{D}\boldsymbol{\sigma}\boldsymbol{D}^T$$
, (14)

and for each row of the resulting D, we retain only the τ entries with the largest absolute values. This process is iterated N times.

Due to noise effects on data, (14) might not have a solution, therefore, some relaxation on the constraint might be required. This is done as follows:

$$\underset{D}{\operatorname{arg\,min}} \{ \sum_{i=1}^{n} (1 - e^{-\frac{s_{i}^{2}}{\sigma^{2}}}) + \lambda \sum_{i=1}^{n} (1 - e^{-\frac{d_{(i,i)}^{2}}{\sigma^{2}}}) \}$$
subject to $\|X - DX\|_{F}^{2} \le \epsilon_{1}$,
$$\|\Sigma - D\sigma D^{T}\|_{F}^{2} \le \epsilon_{2}$$
,
$$(15)$$

where ϵ_1 and ϵ_2 can be tuned to result in the best result.

Solving (13) requires an initial estimate for D, and the final value of the objective function depends on this initial estimate. To obtain the optimal solution, we propose executing the algorithm multiple times, each with a distinct random initialization. The solution that yields the lowest value of the objective function is then retained as the final result. The SLCD algorithm pseudocode is presented in appendix A.

5 SLCD, BEYOND LINEARITY

The proposed framework can be further extended to scenarios in which the SCMs governing the causal relations are nonlinear. The idea relies on the Taylor series expansion of the governing function. Suppose $x_i = h(x)$, where x_i is the ith entry of x and is causally related to x through the deterministic function $h: \mathbb{R}^n \to \mathbb{R}$. Assuming $h \in C^{\infty}(\mathbb{R}^n)$ (the set of infinitely differentiable functions on \mathbb{R}^n), the Taylor series expansion implies that x_i can be expressed as a polynomial in the entries of x. Based on this observation, the following theorem establishes the induced covariance for this scenario.

Theorem 2. Suppose x = g(x), where $x \in \mathbb{R}^n$ is a vector of random variables and $g : \mathbb{R}^n \to \mathbb{R}^n$ represents the causal relations such that $g_i \in C^{\infty}(\mathbb{R}^n)$ for all $i \in \{1, 2, \dots, n\}$. Then, x can be represented as $x = \sum_{i=1}^{\infty} D_i x^i$, where $D_i \in \mathbb{R}^{n \times n}$ denotes the coefficient matrices for all $i \in \{1, 2, \dots, n\}$, and $x^i = (x_1^i, \dots, x_n^i)^{\top}$ is the vector obtained by raising each entry of x to the i-th power. The covariance matrix of x is then given by

$$\Sigma = \sum_{i=1}^{\infty} \sum_{j=i}^{\infty} D_i \sigma_{ij} D_j^T, \tag{16}$$

where σ_{ij} is a diagonal matrix with its l^{th} diagonal elements be $\mathbb{E}[x_l^i x_l^j]$ for all $l \in \{1, 2, \dots, n\}$.

The proof is deferred to Appendix C.1. Similar to the linear case discussed previously, Theorem 2 provides a way to formulate an optimization problem with a constraint stronger than simply minimizing the reconstruction error, i.e., regression. To formulate the optimization problem for the nonlinear case, let $X \in \mathbb{R}^{n \times m}$ be the dataset containing m samples. One can then determine the coefficient matrices that satisfy the following equations:

$$X = \sum_{i=1}^{\infty} D_i X^i, \tag{17}$$

$$\Sigma = \sum_{i=1}^{\infty} \sum_{i=j}^{\infty} D_i \sigma_{ij} D_j^T,$$
(18)

where X^i denotes the elementwise i-th power of X, Σ is the covariance matrix of the data, and σ_{ij} is a diagonal matrix whose l^{th} diagonal entry is given by $\mathbb{E}[x_l^i x_l^j]$ for all $l \in \{1, 2, \dots, n\}$.

6 INDUCED-COVARIANCE: THEORETICAL ANALYSIS

In this section, a theoretical analysis of the proposed method is presented by examining the behavior of solutions that satisfy both the reconstruction and induced covariance constraints. The analysis focuses on the local properties of these solutions, in particular on whether they correspond to isolated points or manifolds, as well as on their sensitivity to variations in the data, with special attention to changes in the covariance matrix.

Theorem 3. Let $F(D) = \Sigma - D\sigma D^T$, where $\Sigma \in \mathbb{R}^{n \times n}$ is the covariance matrix of the data and $\sigma \in \mathbb{R}^{n \times n}$ is the diagonal matrix whose entries are the diagonal elements of Σ . Suppose D^* is a solution of F(D) = 0 that satisfies $D^*X = X$, where $X \in \mathbb{R}^{n \times m}$ denotes the data matrix. Define $\mathbb{S} = \{\Delta \in \mathbb{R}^{n \times n} \mid \Delta X = 0\}$. If there exists no $\Delta \in \mathbb{S} \setminus \{0\}$ such that $D^*\sigma \Delta^T$ is skew-symmetric, then there exists a radius r > 0 such that, for any solution D of F(D) = 0 with DX = X and $\|D^* - D\| < r$, it necessarily follows that $D = D^*$.

The proof is provided in Appendix C.2. This theorem establishes that, under suitable conditions on $D^*\sigma\Delta^T$ along the tangent directions of the feasible solutions (i.e., those satisfying $\Delta X=0$), the solution of F(D)=0 is locally unique. In other words, there exist no other solutions of F(D)=0 satisfying DX=X within a ball of radius r centered at D^* . Next, the following theorem analyzes the sensitivity of the solutions of F(D)=0 to perturbations of the data covariance matrix.

Theorem 4. Let $F(D) = \Sigma - D\sigma D^T$, where $\Sigma \in \mathbb{R}^{n \times n}$ is the covariance matrix of the data and $\sigma \in \mathbb{R}^{n \times n}$ is the diagonal matrix whose entries are the diagonal elements of Σ . Suppose D^* is a solution of F(D) = 0 satisfying $D^*X = X$, where $X \in \mathbb{R}^{n \times m}$ denotes the data matrix. Define $\mathbb{S} = \{\Delta \in \mathbb{R}^{n \times n} \mid \Delta X = 0\}$. If there exists no $\Delta \in \mathbb{S} \setminus \{0\}$ such that $D^*\sigma\Delta^T$ is skew-symmetric, then consider D as a solution of $\Sigma + \Delta \Sigma - D\sigma D^T = 0$, where $\Delta \Sigma$ is a perturbation in the covariance matrix of the data. If $\|\Delta \Sigma\|_F < \frac{c^2}{4\|\sigma\|_2}$, then the following inequality cannot hold:

$$\frac{c - \sqrt{c^2 - 4\|\boldsymbol{\sigma}\|_2 \|\Delta \boldsymbol{\Sigma}\|_F}}{2\|\boldsymbol{\sigma}\|_2} < \|\boldsymbol{D} - \boldsymbol{D}^*\|_F < \frac{c + \sqrt{c^2 - 4\|\boldsymbol{\sigma}\|_2 \|\Delta \boldsymbol{\Sigma}\|_F}}{2\|\boldsymbol{\sigma}\|_2},$$
(19)

where c is the smallest singular value of $D^*\sigma\Delta^T + \Delta\sigma D^{*T}$ on $\mathbb S$.

The proof is provided in Appendix C.3. This theorem establishes a non-feasible region for the distance between the perturbed and unperturbed solutions of the induced covariance equation.

If there exist values of $\Delta \in \mathbb{S} \setminus \{0\}$ such that $D^*\sigma\Delta^T$ is skew-symmetric, then the nearby solutions may not be isolated; that is, there may exist manifolds of solutions along those directions. This necessitates additional restrictions on the solutions to distinguish these directions, making the use of regularization essential in such cases. For this reason, SLCD penalizes the solutions based on their rank and trace in order to further enforce constraints that yield isolated solutions.

7 SIMULATION RESULTS AND ANALYSIS

This section presents the results of our simulation studies. We compare our method with PC, GES, LINGAM IC, LINGAM Direct and BIC exact search for performance evaluation. For comprehen-

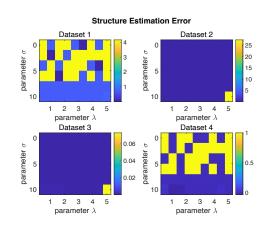


Figure 1: Structure estimation error for various datasets and various hyperparameters.

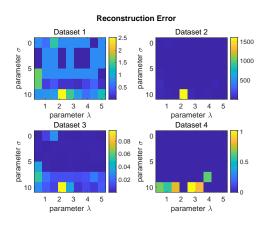


Figure 2: Reconstruction error for various datasets and various hyperparameters.

sive reporting, we evaluate following metrics: the data reconstruction error, the recovery error of the causal matrix, the recovery error of the covariance matrix, precision (porportion of the number of correct estimated links to the total number of estimated links), and recall (porportion of the number of correct estimated links to the total number of links in the true graph). Let $\hat{D} \in \mathbb{R}^{n \times n}$ be the estimated matrix obtained from the proposed algorithm, and $X \in \mathbb{R}^{n \times m}$ represent the training data. The reconstruction error is then defined as:

$$\frac{1}{nm} \|\boldsymbol{X} - \hat{\boldsymbol{D}}\boldsymbol{X}\|_F^2. \tag{20}$$

We define the true structural matrix as $D \in \mathbb{R}^{n \times n}$ and, thus, the recovery error of structural matrix will be:

$$\frac{1}{n^2} \| \mathbf{D} - \hat{\mathbf{D}} \|_F. \tag{21}$$

Let $\Sigma \in \mathbb{R}^{n \times n}$ represent the true covariance matrix of the original data with $\Sigma \in \mathbb{R}^{n \times n}$, then the recovery error of the covariance matrix is defined as

$$\frac{1}{n^2} \| \mathbf{\Sigma} - \hat{\mathbf{D}} \sigma \hat{\mathbf{D}}^T \|_F, \tag{22}$$

in which, σ is a diagonal matrix with diagonal elements of Σ .

Dataset	IV count ¹	x_1	x_2	x_3	Other Entries
Dataset 1	1	U(-2.5, 2.5)	$2x_1$	$0.4x_1$	-
Dataset 2	2	U(-2.5, 2.5)	U(-2.5, 2.5)	$0.3x_1$	$[x_4] = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
Dataset 3	2	U(-2.5, 2.5)	U(-2.5, 2.5)	$x_1 + 3x_2$	$ \begin{bmatrix} x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 2 & 3 & 0 \\ 0 & 2 & 0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} $
Dataset 4	3	U(-2.5, 2.5)	U(-2.5, 2.5)	N(0,4)	$ \begin{bmatrix} x_4 \\ x_5 \\ x_6 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0.3 \\ 2 & 3 & 0 \\ 0 & 2 & 0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} $
Dataset 5	3	U(-2.5, 2.5)	U(-2.5, 2.5)	N(0,4)	$ \begin{bmatrix} x_4 \\ x_5 \\ x_6 \\ x_7 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0.5 \\ 0 & 1 & 2 \\ 1 & 0 & 3 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} $

Table 1: Datasets Information.

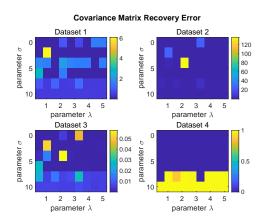


Figure 3: Covariance matrix estimation error for various datasets and various hyperparameters.

For simulation purposes, five distinct datasets were generated, henceforth referred to as Dataset 1, Dataset 2, Dataset 3, Dataset 4, and Dataset 5. Each dataset comprises 1000 samples. Table 1 provides detailed information on the generation process for each dataset. The variables x_i , where $i \in \{1, 2, \ldots, 7\}$, represent the elements of the data vector $\mathbf{x} = [x_1, x_2, \ldots, x_7]^T$. The presence of a '-' symbol in place of a variable indicates its absence from the corresponding dataset, reflecting the varying dimensionality across datasets. The table indicates the data distribution from which the samples of independent variables are drawn. For dependent variables, the table specifies the linear combinations used to generate them.

U(a,b) represents the uniform distribution of data in the [a,b] interval. $N(\mu,\sigma^2)$ represents a Gaussian random variable with mean μ and variance σ^2 . By constructing the datasets in this manner, the variables exhibit the linear sparse relations that SLCD is specifically designed to handle. This approach also enables the evaluation of algorithm performance across various data dimensions. Additionally, the datasets includes independent variables with different data distributions, allowing for the assessment of algorithm robustness under diverse distributional scenarios. It is important to note that for the dependent variables, each linear combination results in a convolution of the data distributions, further contributing to the variability in the distributions of the dataset's variables.

Figures 1 through 3 display the algorithm's simulation results on performance metrics for various hyperparameter settings. The results reveal moderate sensitivity to hyperparameters, with effective recovery of the underlying structure when parameters are chosen appropriately. The figures also indicate a broad range of satisfactory parameters, demonstrating the method's robustness. The detailed performance of SLCD in recovering the structural matrix of each dataset for the hyperparameter pair

#	PC	GES	LG IC	LG Direct	BIC Search	SLCD			
Precision	0.33	0.5	0	0.33	0	0			
Recall	1	0.5	0	0.5	0	0			
Number of Correct link estimation	2	1	0	1	0	0			
Dataset 1									
#	PC	GES	LG IC	LG Direct	BIC Search	SLCD			
Precision	0.5	0.6	0.25	0	0.75	1			
Recall	0.66	1	0.33	0	1	1			
Number of Correct link estimation	2	3	1	0	3	3			
Dataset 2									
#	PC	GES	LG IC	LG Direct	BIC Search	SLCD			
Precision	0.37	0.43	0	0	0.43	1			
Recall	0.6	0.6	0	0	0.6	1			
Number of Correct link estimation	3	3	0	0	3	5			
Dataset 3									
#	PC	GES	LG IC	LG Direct	BIC Search	SLCD			
Precision	1	1	0.2	0.1	0.67	1			
Recall	1	1	0.33	0.17	1	1			
Number of Correct link estimation	6	6	2	1	6	6			
Dataset 4									
#	PC	GES	LG IC	LG Direct	BIC Search	SLCD			
Precision	0.3	0.75	0.08	0.13	0.54	1			
Recall	0.37	0.75	0.12	0.25	1	1			
Number of Correct link estimation	3	6	1	2	6	8			
Dataset 5									

Table 2: Performance comparison of PC, GES, LG IC (LINGAM IC), LG Direct (LINGAM Direct), BIC Exact Search, and SLCD algorithms.

 $(\sigma, \lambda) = (0.3, 5)$ is presented in Table 3 in the Appendix B. It shows that the method successfully recovers the structural matrix of all datasets, with the exception of Dataset 1.

Table 2 presents the simulation results of SLCD in comparison with several well-known causal discovery algorithms. The results indicate that SLCD outperforms the other methods by an average of 35% in precision and 41.5% in recall across Datasets 2 through 5. While all methods exhibit challenges with Dataset 1, SLCD consistently demonstrates superior performance in the remaining datasets.

SLCD demonstrate suboptimal performance on Dataset 1. This can be attributed to the structure of Dataset 1, wherein only one independent variable exists, and all other variables are scalar multiples thereof. This configuration does not provide sufficient information to unambiguously identify the independent variable, as any of the variables could potentially fulfill this role. This ambiguity introduces uncertainty into the algorithm, potentially leading to diverse solutions. However, as the structural complexity increases with the introduction of additional independent variables, the informational content of the data becomes more robust, facilitating more accurate recovery of the underlying causal structure.

8 Conclusion

This paper proposes an algorithm for causal discovery within a linear sparse structure, leveraging properties of the causal structure matrix, specifically its rank, which reflects the number of independent variables, and the notion of induced covariance. Simulation studies confirm the algorithm's effectiveness across diverse configurations. Our next direction is to refine the independence criteria. While induced covariance was useful, it does not fully ensure independence across all the scenarios. Addition of stronger constraints can further enhance this framework. This extension would enhance the algorithm's versatility and reliability across a wider range of applications and data types.

THE USE OF LARGE LANGUAGE MODELS (LLMS) STATEMENT

Large language models were used exclusively for language refinement, including proofreading and grammatical correction.

490 REPRODUCIBILITY STATEMENT

To ensure reproducibility, we provide the pseudocode of the proposed algorithm in Appendix A. The complete implementation code and the datasets used in our experiments are available in the supplementary material. In addition, proofs of the theorems are presented either in the main text or in Appendix C.

REFERENCES

- Gustau Camps-Valls, Andreas Gerhardus, Urmi Ninad, Gherardo Varando, Georg Martius, Emili Balaguer-Ballester, Ricardo Vinuesa, Emiliano Diaz, Laure Zanna, and Jakob Runge. Discovering causal relations and equations from data. *Physics Reports*, 1044:1–68, 2023.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- Christopher Meek. *Graphical Models: Selecting causal and statistical models*. PhD thesis, Carnegie Mellon University, 1997.
- Hosein Mohimani, Massoud Babaie-Zadeh, and Christian Jutten. A fast approach for overcomplete sparse decomposition based on smoothed ℓ^0 norm. *IEEE Transactions on Signal Processing*, 57 (1):289–301, 2009. doi: 10.1109/TSP.2008.2007606.
- Judea Pearl and Thomas S Verma. A theory of inferred causation. In *Studies in Logic and the Foundations of Mathematics*, volume 134, pp. 789–811. Elsevier, 1995.
- Bernhard Schölkopf and Julius von Kügelgen. From statistical to causal learning. In *Proceedings of the International Congress of Mathematicians*, pp. 1, 2022.
- Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- Joni Shaska and Urbashi Mitra. Causal link discovery with unequal edge error tolerance. *IEEE Transactions on Signal Processing*, 73:2848–2861, 2025. doi: 10.1109/TSP.2025.3585825.
- Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10), 2006.
- Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara, Takashi Washio, Patrik O Hoyer, Kenneth Bollen, and Patrik Hoyer. Directlingam: A direct method for learning a linear non-gaussian structural equation model. *Journal of Machine Learning Research-JMLR*, 12(Apr):1225–1248, 2011.
- Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press, 2001.
- Burak Varıcı, Emre Acartürk, Karthikeyan Shanmugam, Abhishek Kumar, and Ali Tajer. Score-based causal representation learning: Linear and general transformations. *arXiv preprint arXiv:2402.00849*, 2024.
- Changhe Yuan and Brandon Malone. Learning optimal bayesian networks: A shortest path perspective. *Journal of Artificial Intelligence Research*, 48:23–65, 2013.

A ALGORITHM PSEUDOCODE

The SLCD pseudocode:

Algorithm 1 Sparse Linear Causal Discovery (SLCD) Algorithm.

```
\boldsymbol{X} \in \mathbb{R}^{n \times m}, N, M, \lambda, \sigma, \tau
\mathbf{for}\ t = 1: M\ \mathbf{do}
      Initialize:
       \mathbf{D}_0 \in \mathbb{R}^{n \times n} : randomly
       if (t == 1) then
              J_{min} \leftarrow J(\boldsymbol{D}_0)
              oldsymbol{D}_{opt} \leftarrow oldsymbol{D}_0
       end if
       \quad \mathbf{for} \ k = 1: N \ \mathbf{do}
              D \leftarrow \text{Solve (14) (e.g. fmincon (MATLAB))}
             if J_{min} > J(\mathbf{D}) then
                     J_{min} \leftarrow J(\boldsymbol{D})
                     oldsymbol{D}_{opt} \leftarrow oldsymbol{D}
             end if
       end for
end for
return oldsymbol{D}_{opt}
```

Dataset	Structure matrix	Estimated structure matrix	(σ, λ)
Dataset 1	$\begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 0.4 & 1 & 0 \end{bmatrix}$	$\begin{bmatrix} 0 & 0.5 & 4.4 \times 10^{-7} \\ 2 & 1 & 1.1 \times 10^{-6} \\ -1.2 \times 10^{-7} & 0.2 & 0 \end{bmatrix}$	(0.3, 5)
Dataset 2	$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0.3 & 1 & 0 & 0 \\ 1 & 2 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & -8.8 \times 10^{-4} & 0 \\ 0 & 1 & -3.3 \times 10^{-4} & 0 \\ 0.3 & 0 & -2.7 \times 10^{-4} & 0 \\ 1.002 & 1.999 & 0 & 0 \end{bmatrix}$	(0.3, 5)
Dataset 3	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 0.999 & 0.0497 & 0 & 0 & 0 \\ 0 & 1.000 & 0 & 0.0102 \\ 0.976 & 3.049 & 0 & 0 & 0 \\ -0.0147 & 1.999 & 0 & 0 & 0 \\ 1.990 & 1.099 & 0 & 0 & 0 \end{bmatrix}$	(0.3, 5)
Dataset 4	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0.3 & 0 & 0 & 0 \\ 2 & 3 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0.5 & 0 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 0.999 & -0.009 & 0 & 0 & 0 & 0 \\ 0.016 & 0.999 & 0 & 0 & 0 & 0 \\0432 & 0 & 0.997 & 0 & 0 & 0 \\ 0.987 & 0 & 0.3019 & 0 & 0 & 0 \\ 2.048 & 2.982 & 0 & 0 & 0 & 0 \\ 0 & 1.995 & 0.483 & 0 & 0 & 0 \end{bmatrix}$	(0.3, 5)
Dataset 5	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0.5 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 0 & 0 & 0 & 0 \\ 1 & 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 0.997 & .0525 & 0 & 0 & 0 & 0 \\ -0.082 & 0.994 & 0 & 0 & 0 & 0 & 0 \\ 0.057 & 0 & 0.998 & 0 & 0 & 0 & 0 \\ 1.025 & 0 & 0.491 & 0 & 0 & 0 & 0 \\ 0 & 0.956 & 2.024 & 0 & 0 & 0 & 0 \\ 1.168 & 0 & 2.986 & 0 & 0 & 0 & 0 \\ 0 & 0.975 & 1.025 & 0 & 0 & 0 & 0 \end{bmatrix}$	(0.3, 5)

Table 3: True structural matrix and the output of SLCD.

B STRUCTURAL RECOVERY RESULTS OF SLCD ACROSS DIFFERENT DATASETS

Table 3 shows the recovered structural matrix usign SLCD as well as the ground truth structural matrix for each dataset with $(\sigma, \lambda) = (0.3, 5)$.

C PROOFS

C.1 Proof of Theorem 2

Let $x \in \mathbb{R}^n$ be a random vector with the following causal structure:

$$x = \sum_{i=1}^{\infty} D_i x^i, \tag{23}$$

where x^i is the vector obtained by elementwise raising of x to the i-th power, and D_i is the corresponding coefficient matrix.

To calculate the covariance matrix of the data, we need to compute $\mathbb{E}[x_i x_j]$. Using equation 23, we have

$$\mathbb{E}[x_i x_j] = \mathbb{E}\left[\left(\sum_{l=1}^{\infty} \boldsymbol{d}_i^{(l)T} \boldsymbol{x}^l\right) \left(\sum_{k=1}^{\infty} \boldsymbol{d}_j^{(k)T} \boldsymbol{x}^k\right)\right],\tag{24}$$

where $m{d}_i^{(l)T}$ and $m{d}_j^{(k)T}$ are the $i^{ ext{th}}$ and $j^{ ext{th}}$ rows of the matrix $m{D}_l$ and $m{D}_k$, respectively.

Applying this procedure to all pairs of entries, we obtain

$$\Sigma = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} D_i \sigma_{ij} D_j^T, \tag{25}$$

where σ_{ij} is a diagonal matrix whose l^{th} diagonal element is $\mathbb{E}[x_i^i x_l^j]$ for all $l \in \{1, 2, \dots, n\}$.

C.2 PROOF OF THEOREM 3

Suppose $F(D) = \Sigma - D\sigma D^T$ and D^* is a solution of F(D) = 0 that satisfies $D^*X = X$. Let D be another solution of F(D) = 0 satisfying DX = X, and define $\Delta = D - D^*$. Then $\Delta X = 0$ and

$$F(\mathbf{D}) = F(\mathbf{D}^* + \mathbf{\Delta}) = \mathbf{\Sigma} - (\mathbf{D}^* + \mathbf{\Delta})\sigma(\mathbf{D}^* + \mathbf{\Delta})^T$$
(26)

$$=\underbrace{\Sigma - D^* \sigma D^{*T}}_{0} - \underbrace{\left(D^* \sigma \Delta^T + \Delta \sigma D^{*T}\right)}_{\mathcal{L}(\Delta)} - \underbrace{\Delta \sigma \Delta^T}_{\mathcal{Q}(\Delta)}$$
(27)

$$= -\mathcal{L}(\Delta) - \mathcal{Q}(\Delta) \tag{28}$$

$$=0, (29)$$

which yields the identity

$$\mathcal{L}(\Delta) + \mathcal{Q}(\Delta) = 0, \quad \Delta X = 0.$$
 (30)

Let $\mathbb{S} = \{ \Delta \in \mathbb{R}^{n \times n} \mid \Delta X = \mathbf{0} \}$. By assumption, for no $\Delta \in \mathbb{S} \setminus \{\mathbf{0}\}$, $D^* \sigma \Delta^T$ is skew-symmetric, which implies that $\mathcal{L}(\Delta) = \mathbf{0}$ on \mathbb{S} , only happens if $\Delta = \mathbf{0}$. This results in the existence of a positive constant c defined as

$$c = \inf_{\substack{\Delta \in \mathbb{S} \\ \|\Delta\|_F = 1}} \|\mathcal{L}(\Delta)\|_F. \tag{31}$$

Hence, for all $\Delta \in \mathbb{S}$,

$$\|\mathcal{L}(\mathbf{\Delta})\|_F > c \|\mathbf{\Delta}\|_F. \tag{32}$$

Using norm inequalities for $\mathcal{Q}(\Delta)$, we have

$$\|\mathcal{Q}(\mathbf{\Delta})\|_F = \|\mathbf{\Delta}\boldsymbol{\sigma}\mathbf{\Delta}^T\|_F \le \|\boldsymbol{\sigma}\|_2 \|\mathbf{\Delta}\|_F^2, \quad \forall \mathbf{\Delta} \in \mathbb{S}.$$
 (33)

Combining these results with equation 30 gives

$$c \|\Delta\|_F \le \|\mathcal{L}(\Delta)\|_F = \|\mathcal{Q}(\Delta)\|_F \le \|\sigma\|_2 \|\Delta\|_F^2,$$
 (34)

which leads to

$$\|\Delta\|_F \ge \frac{c}{\|\sigma\|_2}.\tag{35}$$

This final inequality establishes a positive lower bound for $\|\Delta\|_F = \|D - D^*\|_F$, indicating that any two feasible solutions must be separated by at least $\frac{c}{\|\sigma\|_2}$, which implies that the solution is isolated and locally unique.

C.3 PROOF OF THEOREM 4

Suppose $F(D) = \Sigma - D\sigma D^T$ and D^* is a solution of F(D) = 0 that satisfies $D^*X = X$. Let the data covariance matrix be perturbed as $\Sigma \to \Sigma + \Delta \Sigma$. Assume D is the solution of F(D) with the perturbed covariance matrix, i.e.,

$$\Sigma + \Delta \Sigma - D \sigma D^T = 0,$$

and satisfies DX = X. Let $\mathbb{S} = \{\Delta \in \mathbb{R}^{n \times n} \mid \Delta X = 0\}$. Then $D - D^* = \Delta \in \mathbb{S}$, and

$$\mathbf{0} = \mathbf{\Sigma} + \Delta \mathbf{\Sigma} - \mathbf{D} \boldsymbol{\sigma} \mathbf{D}^T \tag{36}$$

$$= \Sigma + \Delta \Sigma - (D^* + \Delta)\sigma(D^* + \Delta)^T$$
(37)

$$=\underbrace{\Sigma - D^* \sigma D^{*T}}_{0} - \underbrace{\left(D^* \sigma \Delta^T + \Delta \sigma D^{*T}\right)}_{\mathcal{L}(\Delta)} - \underbrace{\Delta \sigma \Delta^T}_{\mathcal{Q}(\Delta)} + \Delta \Sigma, \tag{38}$$

which gives

$$\mathcal{L}(\Delta) + \mathcal{Q}(\Delta) = \Delta \Sigma, \quad \Delta X = 0.$$
 (39)

Following the procedure in C.2, the nonexistence of any $\Delta \in \mathbb{S} \setminus \{0\}$ making $D^*\sigma\Delta^T$ skew-symmetric implies the following inequalities

$$\|\mathcal{L}(\Delta)\|_F \ge c \|\Delta\|_F, \quad \forall \Delta \in \mathbb{S},$$
 (40)

$$\|\mathcal{Q}(\boldsymbol{\Delta})\|_F \le \|\boldsymbol{\sigma}\|_2 \|\boldsymbol{\Delta}\|_F^2, \quad \forall \boldsymbol{\Delta} \in \mathbb{S},$$
 (41)

where

$$c = \inf_{\substack{\Delta \in \mathbb{S} \\ \|\Delta\|_F = 1}} \|\mathcal{L}(\Delta)\|_F. \tag{42}$$

Then we have

$$\mathcal{L}(\Delta) = \Delta \Sigma - \mathcal{Q}(\Delta), \quad \Delta X = 0, \tag{43}$$

$$\|\mathcal{L}(\Delta)\|_F = \|\Delta\Sigma - \mathcal{Q}(\Delta)\|_F,\tag{44}$$

$$\|\mathcal{L}(\mathbf{\Delta})\|_F \le \|\Delta \mathbf{\Sigma}\|_F + \|\mathcal{Q}(\mathbf{\Delta})\|_F,\tag{45}$$

$$c \|\mathbf{\Delta}\|_F \le \|\mathcal{L}(\mathbf{\Delta})\|_F \le \|\Delta\mathbf{\Sigma}\|_F + \|\mathcal{Q}(\mathbf{\Delta})\|_F \le \|\Delta\mathbf{\Sigma}\|_F + \|\boldsymbol{\sigma}\|_2 \|\mathbf{\Delta}\|_F^2, \tag{46}$$

$$c \|\mathbf{\Delta}\|_F \le \|\Delta\mathbf{\Sigma}\|_F + \|\boldsymbol{\sigma}\|_2 \|\mathbf{\Delta}\|_F^2. \tag{47}$$

Let $r = \|\Delta\|_F$. Then the quadratic inequality

$$\|\boldsymbol{\sigma}\|_2 r^2 - cr + \|\Delta \boldsymbol{\Sigma}\|_F > 0 \tag{48}$$

holds. If the discriminant of equation 48, $c^2-4\|\pmb{\sigma}\|_2\|\Delta\pmb{\Sigma}\|_F$, is less than zero, the inequality has no real roots and imposes no bounds on r. However, if $\|\Delta\pmb{\Sigma}\|_F<\frac{c^2}{4\|\pmb{\sigma}\|_2}$, then equation 48 implies that r cannot lie between the two real roots, i.e.,

$$\frac{c - \sqrt{c^2 - 4\|\boldsymbol{\sigma}\|_2 \|\Delta \boldsymbol{\Sigma}\|_F}}{2\|\boldsymbol{\sigma}\|_2} < r < \frac{c + \sqrt{c^2 - 4\|\boldsymbol{\sigma}\|_2 \|\Delta \boldsymbol{\Sigma}\|_F}}{2\|\boldsymbol{\sigma}\|_2}$$
(49)

is infeasible.