

# THE ROBUSTNESS OF DIFFERENTIABLE CAUSAL DISCOVERY IN MISSPECIFIED SCENARIOS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Causal discovery aims to learn causal relationships between variables from targeted data, making it a fundamental task in machine learning. However, causal discovery algorithms often rely on unverifiable causal assumptions, which are usually difficult to satisfy in real-world data, thereby limiting the broad application of causal discovery in practical scenarios. Inspired by these considerations, this work extensively benchmarks the empirical performance of various mainstream causal discovery algorithms, which assume i.i.d. data, under eight model assumption violations. Our experimental results show that differentiable causal discovery methods exhibit robustness under the metrics of Structural Hamming Distance and Structural Intervention Distance of the inferred graphs in **commonly used** challenging scenarios, except for scale variation. We also provide the theoretical explanations for the performance of differentiable causal discovery methods. Finally, our work aims to comprehensively benchmark the performance of recent differentiable causal discovery methods under model assumption violations, and provide the standard for reasonable evaluation of causal discovery, as well as to further promote its application in real-world scenarios.

## 1 INTRODUCTION

In the realm of modern science, numerous endeavors hinge upon the elucidation of underlying causal mechanisms. However, owing to practical constraints including costs, risks, and ethical implications, the execution of randomized experiments frequently proves unviable. Consequently, mining causal relationships from purely observational data, known as causal discovery, plays a crucial role in addressing causal questions such as intervention and counterfactual (Peters et al., 2017; Spirtes et al., 2001; Pearl, 2009; Pearl & Mackenzie, 2018).

Causal discovery encompasses a comprehensive suite of methodologies, primarily categorized into constraint-based, score-based, functional causal model-driven, and gradient-based approaches. These methods often rely on unverifiable causal assumptions as their foundation (Peters et al., 2017; Vowels et al., 2022). Constraint-based methods, notably PC (Spirtes & Glymour, 1991) and FCI (Spirtes et al., 1995), meticulously reconstruct causal graphs to the Markov equivalence class (MEC) through rigorous statistical independence tests, guided by the faithfulness assumption. Score-based techniques, such as GES (Chickering, 2002), employ a scoring function to quantify the congruence between an equivalence class and observed data, optimally searching the vast landscape of directed acyclic graphs (DAGs) to identify the MEC.

To transcend the limitation of solely identifying MECs from observational data, functional causal model-based methods, exemplified by LiNGAM (Shimizu et al., 2006), leverage precise assumptions regarding the functional class and noise distribution within structural equation models (SEMs), enabling the unambiguous identification of a unique DAG. **Recently, Zheng et al. (2018) introduced gradient-based techniques (e.g. NOTEARS (Zheng et al., 2018)), which convert combinatorial acyclic constraints into smooth equality constraints and solve the optimization by transforming equality-constrained optimization into unconstrained optimization through the augmented Lagrangian method (ALM) (Nemirovsky, 1999). In some literature (Zhang et al., 2023; Liu et al., 2023), gradient-based methods are also referred to as differentiable causal discovery.**

Apart from the various assumptions of the methods above, traditional approaches typically rely on causal sufficiency and no measurement error assumptions to simplify the problem (Peters et al., 2017;

054 [Zhang et al., 2018](#)). Real-world data often fail to meet all of these assumptions, and these are also  
055 impossible to verify adequately (Peters et al., 2017). Although some studies have considered the  
056 complexity of real data and developed algorithms targeted at latent confounders (Spirtes et al., 1995;  
057 Xie et al., 2020; Salehkaleybar et al., 2020; Cai et al., 2019; Kong et al., 2023; Cai et al., 2023),  
058 measurement error (Zhang et al., 2018; Dai et al., 2022), heterogeneity (Huang et al., 2020; Cai et al.,  
059 2020; Ghassami et al., 2017; 2018), [scale variation](#) (Shimizu et al., 2011; Reisach et al., 2023; Deng  
060 et al., 2024), and missing data (Tu et al., 2019a; Gao et al., 2022), the true mechanisms remain unclear  
061 when [the causal discovery algorithms](#) applied to real data. These specifically designed algorithms  
062 also cannot be effectively employed for real data. Therefore, the robustness of causal discovery  
063 algorithms in scenarios where model hypotheses are violated is of great importance.

064 Previous research (Heinze-Deml et al., 2018) mainly evaluated various constraint-based and score-  
065 based algorithms under different scenarios, only focusing on linear SEM. The work (Mooij et al.,  
066 2016) benchmarked causal discovery for nonlinear additive noise models and information-geometric  
067 approaches, limiting to bivariate scenarios. The previous study (Singh et al., 2017) primarily assessed  
068 algorithms that use only observational data, a mix of observational and interventional data, and active  
069 learning, but their algorithm outputs were restricted to MEC. Also, those works (Glymour et al., 2019;  
070 Vowels et al., 2022) reviewed the advancements in traditional causal discovery (constraint-based,  
071 score-based, and functional causal model-based) and differentiable causal discovery, respectively, but  
072 lacked experimental support. The recent work (Ng et al., 2024) conducted an experimental assessment  
073 of the advancements in differentiable causal discovery, illuminating the shortcomings of current  
074 methods. However, their evaluation overlooked the ubiquitous violation of model assumptions that  
075 characterize real-world applications. Conversely, another work (Montagna et al., 2023) benchmarked  
076 the efficacy of traditional causal discovery algorithms, encompassing score-matching techniques,  
077 under scenarios where model assumptions were violated. Nevertheless, their analysis did not  
078 encompass the recent strides made in differentiable causal discovery, and the misspecified conditions  
079 they evaluated were constrained in scope. Given that the application of causal discovery methods to  
080 real data inevitably entails the violation of one or more unidentified assumptions, and that algorithms  
081 premised on specific assumptions may falter in practical use, the robustness of causal discovery in  
082 such misspecified contexts assumes paramount importance.

082 Our study undertakes an exhaustive empirical evaluation of both established and cutting-edge causal  
083 discovery methodologies, comprehensively examining their performance under diverse scenarios  
084 with assumption violations. The misspecified scenarios encompassed in our analysis represent the  
085 most extensive coverage in the current research landscape. We meticulously evaluate mainstream  
086 causal discovery approaches, spanning constraint-based, score-based, functional causal model-based,  
087 gradient-based methodologies, among others, ensuring a holistic view of the field. Notably, our  
088 work fills a crucial gap in the literature by being the first to assess the performance of gradient-based  
089 methods across a wide array of misspecified scenarios. [Considering their practical implementation  
090 potential, it is important to evaluate their performance.](#) Our contributions can be summarized as  
091 follows:

- 092
- 093 • We conduct extensive large-scale experimental evaluations of twelve prominent causal discovery  
094 algorithms across eight pivotal model assumption violation scenarios. Our rigorous research  
095 endeavor involves executing over 70,000 experiments on more than 2,400 synthetic datasets,  
096 ensuring a comprehensive assessment of the algorithm capabilities.
- 097
- 098 • We delve into challenging scenarios such as heterogeneity, scale variation, missing data, and  
099 mechanism violation, thereby enriching the benchmark data landscape for model assumption  
100 violations. This also aims to foster more comprehensive benchmark testing and foster a more  
101 rational evaluation framework for future causal discovery algorithms.
- 102
- 103
- 104 • Our analysis of the experimental outcomes offers theoretical insights into the performance of linear  
105 differentiable causal discovery methods under certain misspecified scenarios. Recognizing the  
106 robustness demonstrated by differentiable methods in these [commonly used](#) challenging settings,  
107 we underscore the significant value of further in-depth research into differentiable causal discovery,  
as it holds the promise of advancing the field in novel and impactful directions.

## 2 BACKGROUND

In this section, we introduce the definition of causal discovery (Section 2.1), functional causal model-based (Section 2.2), score-based (Section 2.3) and gradient-based (Section 2.3) methods. For constraint-based methods, see Appendix A.2 and B.1.

### 2.1 TASK FORMULATION

A structural causal model  $\mathcal{M}$  (Pearl, 2009) consists of the set of endogenous variables  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$ , exogenous variables  $U = (U_1, \dots, U_d) \in \mathbb{R}^d$ , and functional mechanisms  $\mathcal{F} = (f_1, \dots, f_d)$ . Each variable  $X_i$  is defined by a structural equation:

$$X_i = f_i(X_{pa(X_i)}, U_i), \forall i = 1, \dots, d, \quad (1)$$

where  $X_i$  is the  $i$ -th node variable,  $pa(X_i)$  denote the parents of  $X_i$ ,  $f_i : \mathbb{R}^{|X_{pa(X_i)}|+1} \rightarrow \mathbb{R}$  is the causal structure function, and  $U = (U_1, \dots, U_d)$  are jointly independent noise variables with covariance matrix  $\Omega = \text{cov}(U) = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ .

The task of causal discovery is to infer a DAG  $\mathcal{G}$  that describes the causal relationships among variables from  $n$  independent and identically distributed (i.i.d.) observational data  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , which are drawn from the joint probability distribution  $P(X)$ .

### 2.2 STRUCTURE IDENTIFIABILITY

To uniquely identify a DAG  $\mathcal{G}$  from purely observational data  $\mathbf{X}$  sampled from  $P(X)$ , we need to make assumptions about the SEM in (1). Considering a set of assumptions  $\mathcal{A}$  on a structural causal model (SCM)  $\mathcal{M}_{\mathcal{A}} = (P(X), \mathcal{G})$ , the graph  $\mathcal{G}$  is identifiable from  $P(X)$  if there is no other SCM  $\tilde{\mathcal{M}}_{\mathcal{A}} = (\tilde{P}(X), \tilde{\mathcal{G}})$  satisfying the same  $\mathcal{A}$  such that  $\tilde{\mathcal{G}} \neq \mathcal{G}$  and  $\tilde{P}(X) = P(X)$ . Existing identifiable causal models include: linear non-Gaussian acyclic models (Shimizu et al., 2006), linear Gaussian models with equal noise variances (Peters & Bühlmann, 2014), post-nonlinear models (Zhang & Hyvarinen, 2012) and nonlinear additive noise models (Hoyer et al., 2008; Peters et al., 2014).

### 2.3 DIFFERENTIABLE SCORE-BASED CAUSAL DISCOVERY

Traditional score-based causal discovery defines a combinatorial optimization problem:

$$\min_{\mathcal{G}} F(\mathcal{G}; \mathbf{X}) = \mathcal{L}_{\text{rec}}(\mathcal{G}; \mathbf{X}) + \lambda \mathcal{L}_{\text{sparse}}(\mathcal{G}) \quad \text{s.t.} \quad \mathcal{G} \in \text{DAG}, \quad (2)$$

where  $F$  is a score function,  $\mathcal{L}_{\text{rec}}(\mathcal{G}; \mathbf{X})$  represents the goodness-of-fit between the estimated DAG  $\mathcal{G}$  and the true DAG,  $\mathcal{L}_{\text{sparse}}(\mathcal{G})$  denotes the sparsity regularization term and  $\lambda$  is a hyperparameter that controls the strength of regularization.

As the number of nodes rises, the total count of possible DAGs expands super-exponentially (Robinson, 1973). Consequently, most conventional score-based approaches utilize local heuristic search techniques, including greedy search (Chickering, 2002; Hauser & Bühlmann, 2012) and hill-climbing (Gámez et al., 2011; Tsamardinos et al., 2006).

In addition to search strategies, the design of score functions is also crucial. Commonly, score functions are classified into two categories: Bayesian-based scores and information-theoretic scores. Bayesian-based scores emphasize goodness-of-fit and enable the integration of prior knowledge, such as the Bayesian Dirichlet equivalent (Heckerman et al., 1995) and the K2 score (Kayaalp & Cooper, 2012). Information-theoretic scores, on the other hand, account for both model goodness-of-fit and complexity, including the Bayesian information criterion (Neath & Cavanaugh, 2012) and the Akaike information criterion (Akaike, 1998).

To overcome the challenges of combinatorial optimization, NOTEARS (Zheng et al., 2018) formulates the DAG structure learning task as:

$$\min_{\mathcal{G}} F(\mathcal{G}; \mathbf{X}) \quad \text{s.t.} \quad h(W(\mathcal{G})) = 0, \quad (3)$$

where  $W(\mathcal{G}) \in \mathbb{R}^{d \times d}$  is a weighted adjacency matrix,  $d$  is the number of nodes,  $h(W(\mathcal{G})) = 0$  is a differentiable equality DAG constraint.

162  $h(W(\mathcal{G})) = 0$  if and only if  $W(\mathcal{G})$  is a DAG. Commonly used DAG constraints include  $h(W(\mathcal{G})) =$   
 163  $\text{Tr}(e^{W \circ W}) - d$  (Zheng et al., 2018),  $h(W(\mathcal{G})) = \text{Tr}[(I + \alpha W \circ W)^d] - d$  ( $\alpha > 0$ ) (Yu et al., 2019)  
 164 and  $h^s(W(\mathcal{G})) = -\log \det(sI - W \circ W) + d \log s$  ( $s > 0$ ) (Bello et al., 2022). Furthermore, we  
 165 can transform the equality-constrained optimization (3) into unconstrained optimization (4) using the  
 166 ALM (Nemirovsky, 1999):

$$\min_{\mathcal{G}} F(\mathcal{G}; \mathbf{X}) + \alpha_t h(W(\mathcal{G})) + \frac{\mu_t}{2} |h(W(\mathcal{G}))|^2, \tag{4}$$

167  
 168 where  $\alpha_t$  and  $u_t$  are the Lagrange multiplier and penalty parameter at the  $t$ -th iteration, respectively.  
 169

### 170 171 172 173 3 EXPERIMENTAL DESIGN

174 In this section, we introduce the generation of synthetic datasets with violated model assumptions,  
 175 the tested causal discovery algorithms, the algorithm hyperparameters, and the evaluation metrics.  
 176

#### 177 178 179 3.1 SYNTHETIC DATASETS

180 Many prevalent causal discovery approaches hinge upon unverifiable assumptions. Our study pri-  
 181 marily scrutinizes the efficacy of these methodologies in circumstances where their underlying  
 182 assumptions are breached. To achieve this, we commence by elucidating the baseline model under  
 183 both linear and nonlinear frameworks, subsequently delving into scenarios where these fundamental  
 184 assumptions fail to hold.  
 185

186 Table 1: Summary of the algorithm assumptions and their corresponding output graph types. The  
 187 content within the cells indicates whether an algorithm supports (✓) or does not support (✗) the  
 188 specific condition in the corresponding row. The table style is adjusted from Montagna et al. (2023).  
 189

	PC	GES	DirectLINGAM	CAM	Var-SortnRegress	R <sup>2</sup> -SortnRegress	NOTEARS	GOLEM	NOTEARS-MLP	GranDAG	NoCurl	DAGMA
Gaussian noise	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓
Non-Gaussian noise	✓	✗	✓	✗	✓	✓	✓	✗	✓	✗	✓	✓
Linear mechanisms	✓	✓	✓	✗	✓	✓	✓	✓	✗	✗	✓	✓
Nonlinear mechanisms	✓	✓	✗	✓	✗	✗	✗	✗	✓	✓	✓	✓
Unfaithful distribution	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Confounding effects	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Measurement errors	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Autoregressive effects	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Heterogeneous effects	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Scale-variant effects	✗	✗	✓	✗	✗	✓	✗	✗	✗	✗	✗	✗
Missing mechanisms	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Output	CPDAG	CPDAG	DAG	DAG	DAG	DAG	DAG	DAG	DAG	DAG	DAG	DAG

192  
 193  
 194  
 195  
 196  
 197  
 198  
 199  
 200  
 201  
 202  
 203  
**Linear vanilla model.** In linear SCM, following the settings of Zheng et al. (2018), coefficients  
 are sampled from  $U(-2, -0.5) \cup U(0.5, 2)$  with additive standard Gaussian noise. We refer to this  
 model as the linear vanilla model, which satisfies both identifiability and the assumptions of most  
 linear benchmark methods (see Table 1).

204  
 205  
 206  
 207  
 208  
 209  
**Nonlinear vanilla model.** In nonlinear settings, following the settings of Zheng et al. (2020), the  
 SEM in equation (1) is generated under the Gaussian process with radial basis function kernel of  
 bandwidth one, where  $f_i$  is additive noise models with  $U_i$  as a standard Gaussian noise. We refer to  
 this model as the nonlinear vanilla model, which satisfies both identifiability and the assumptions of  
 nonlinear benchmark methods (see Table 1).

210  
 211  
 212  
 213  
 214  
 215  
 To eliminate the impact of Gaussian noise in the vanilla model on experimental results, we also  
 consider cases where the vanilla model uses non-Gaussian noise (see Appendix G).

### 3.1.1 MODEL ASSUMPTION VIOLATION SCENARIOS

Four scenarios of model assumption violations are defined below. The other four cases, i.e., confounded, measurement error, unfaithful and autoregressive model, follow the same settings as Montagna et al. (2023). These eight misspecified scenarios can be applied to both the linear vanilla and nonlinear vanilla model to generate datasets.

**Heterogeneous model.** Existing causal discovery algorithms typically rely on the assumption of i.i.d. data. However, real data often exhibit distribution shifts (Huang et al., 2020). The heterogeneous multi-domain data considered in this paper primarily refers to scenarios where the underlying causal generation process remains unchanged, but the distribution of noise terms varies (Huang et al., 2020; Zhang et al., 2023; Wang et al., 2022). Specifically, we consider data from two domains  $e_1$  and  $e_2$ . The proportion of data from  $e_1$  is  $P_1 \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$ , and the proportion from  $e_2$  is  $1 - P_1$ . The noise variance in  $e_1$  is set the same as the vanilla model, while variance in  $e_2$  is set to  $\gamma \in \{0.01, 0.05, 0.1, 0.5\}$ .

**Scale-variant model.** Reisach et al. (2021) observed a significant performance decline in linear gradient-based methods, such as NOTEARS (Zheng et al., 2018) and GOLEM (Ng et al., 2020), when applied to data with scale variation. However, there has been a notable absence of research investigating the performance of nonlinear methods in the context of scale variation. Therefore, we also consider scale variation as a misspecified scenario. The structural equations considered are:

$$\bar{X}_i = \frac{X_i - u_i}{\sqrt{\text{Var}(X_i)}}, \forall i = 1, \dots, d, \quad (5)$$

where  $u_i$  and  $\text{Var}(X_i)$  are the mean and variance of  $X_i$ , respectively. The input data are standardized, while the ground truth graph remains consistent with the causal graph that generates the original data.

**Missing model.** Missing data is a prevalent challenge in real-world datasets, necessitating that causal discovery algorithms effectively address this issue (Tu et al., 2019b). In our study, we adopt the Missing Completely At Random (MCAR) mechanism (Tu et al., 2019a), where the occurrence of missing values follows a Bernoulli distribution with a missingness probability of  $\beta \in \{0.005, 0.01, 0.05, 0.1\}$ . Given that the algorithms under consideration are incapable of directly processing datasets with missing values, we eliminate any records containing such gaps. To mitigate the influence of data quantity on the experimental outcomes, we ensure the volume of data remains consistent before and after the removal of incomplete records.

**Mechanism violation.** Most current causal discovery algorithms presuppose either linear or nonlinear mechanism, especially methods based on functional causal models (Shimizu et al., 2006; Peters & Bühlmann, 2014; Zhang & Hyvarinen, 2012; Hoyer et al., 2008; Peters et al., 2014). These methods necessitate specific assumptions about the SEM mechanism to guarantee identifiability. Given that the SEM in real-world data is typically unknown, the robustness of algorithms in the face of mechanism violation becomes critically important. In mechanism violation, the input data for algorithms designed for linear SEM will adhere to the nonlinear vanilla model, whereas the input data for algorithms tailored to nonlinear SEM will conform to the linear vanilla model.

### 3.1.2 DATA GENERATION

Following the data generation of Zheng et al. (2018; 2020) and Liu et al. (2023), different datasets are generated for both linear and nonlinear vanilla model. We simulate ER and SF graphs based on the number of nodes  $d \in \{10, 20, 50\}$ , average degree of nodes  $k \in \{2, 4\}$ . In addition, we consider scenarios with Gaussian Random Partitions (GRP) (Brandes et al., 2003) graph and an average node degree of 6. For each experimental configuration and scenario, 10 datasets of 2000 samples are generated. The mean and standard deviations of the evaluation metrics (Section 3.4) is reported to ensure a fair comparison.

## 3.2 METHODS

We select 12 mainstream causal discovery algorithms, including constraint-based, score-based, functional causal model-based, gradient-based and other representative methods. For a more detailed introduction to the various methods, see the Appendix B.

### 3.3 HYPERPARAMETER SETTINGS

PC (Spirtes & Glymour, 1991), CAM (Bühlmann et al., 2014), and GraN-DAG require adjustment of the significance level  $\alpha$  in the statistical independence tests. NOTEARS, GOLEM, NOTEARS-MLP, and DAGMA need adjustment of the sparsity coefficient  $\lambda_1$  for the  $l_1$ -norm regularization term. Typically, the ground truth of real data is unknown, making it difficult to effectively select hyperparameters for various algorithms. Thus, to ensure a fair comparison of various methods, we tune  $\lambda_1$  in  $\{0.005, 0.01, 0.05, 0.5, 2, 5\}$  and tune  $\alpha$  in  $\{0.001, 0.01, 0.05, 0.1\}$ .

### 3.4 EVALUATION METRICS

We employ Structural Hamming Distance (SHD) and Structural Intervention Distance (SID) (Peters & Bühlmann, 2015) to evaluate performance. SHD counts the number of edge insertions, deletions, and reversals necessary to transform the estimated graph into the true graph. SID is used to assess the distinctions in intervention distribution between the estimated and the true graph. Intuitively, SHD focuses on differences in graph structure, while SID focuses on differences in causal ordering. Lower SHD and SID values indicate better estimation of the target causal graph by the algorithm. For cases where the output is a MEC, we follow the same approach as Zheng et al. (2018), evaluating them favorably by assuming the undirected edges in the MEC are in the correct direction.

## 4 CRITICAL EXPERIMENTAL RESULTS AND INSIGHTS

In this section, we first present the experimental results of the misspecified datasets generated according to Section 3.1.1, comparing them with the findings from the vanilla scenario to draw conclusions. Finally, we provide a more in-depth discussion on the performance of CAM (Section 4.1.1) and offer theoretical insights into the performance of differentiable causal discovery (Section 4.1.2). Due to space limitations, the main text focuses on the experimental results for **ER-2 graphs of 10 nodes** (Table 2.1, 2.2, 3.1, 3.2, 4.1, 4.2), whereas similar conclusions apply to **different nodes, graph types and graph densities** (see Appendix E and J). To visually and concisely present the results, the outcomes of the 10 nodes graph under linear, nonlinear, and MLP settings (Section 4.1.1) are summarized in Figure 1. **We also consider the real-world Sachs (Sachs et al., 2005) dataset (see Appendix I), combined misspecified scenarios (see Appendix F), vanilla model with non-Gaussian noise (see Appendix G), semi-synthetic data (see Appendix L) and runtime of the benchmark methods (see Appendix D).** For each scenario, we generate datasets with 10 different random seeds, each time drawing 2000 samples. **We report the mean and standard deviations of the metrics over 10 trials.** To guarantee a fair comparison of various methods, the hyperparameters for each method are determined as the optimal values relative to the specific dataset.

### 4.1 CURRENT METHODS’ PERFORMANCE IN MISSPECIFIED SCENARIOS

Our experiments show that differentiable causal discovery algorithms almost always achieve optimal or competitive performance in **commonly used misspecified** scenarios other than scale variation. In this paper, robustness refers to the ability of the model to perform well in misspecified scenarios, consistent with the understanding of Montagna et al. (2023).

**Confounded, measurement error, autoregressive and heterogeneous model.** Table 2.1, 2.2, 3.1 and 3.2 indicate that under the commonly used **confounded, measurement error, autoregressive and heterogeneous** ( $P_1 = 0.5, \gamma = 0.1$ ) scenarios, differentiable causal discovery achieves optimal or competitive performance compared to other methods. For the nonlinear Gaussian process mechanism, although CAM (Bühlmann et al., 2014) demonstrates better performance, the discussion in Section 4.1.1 reveals that CAM still has limitations compared to differentiable causal discovery.

**Missing model.** We generate missing data that are MCAR with the missingness probability  $\alpha = 0.01$ . Table 2.1, 2.2, 3.1 and 3.2 indicate that under MCAR, the result of various algorithms is close to that in the vanilla model. Our experiments show that the performance of differentiable causal discovery under missing data is consistent with traditional methods considered by Tu et al. (2019b;a), including PC and GES.

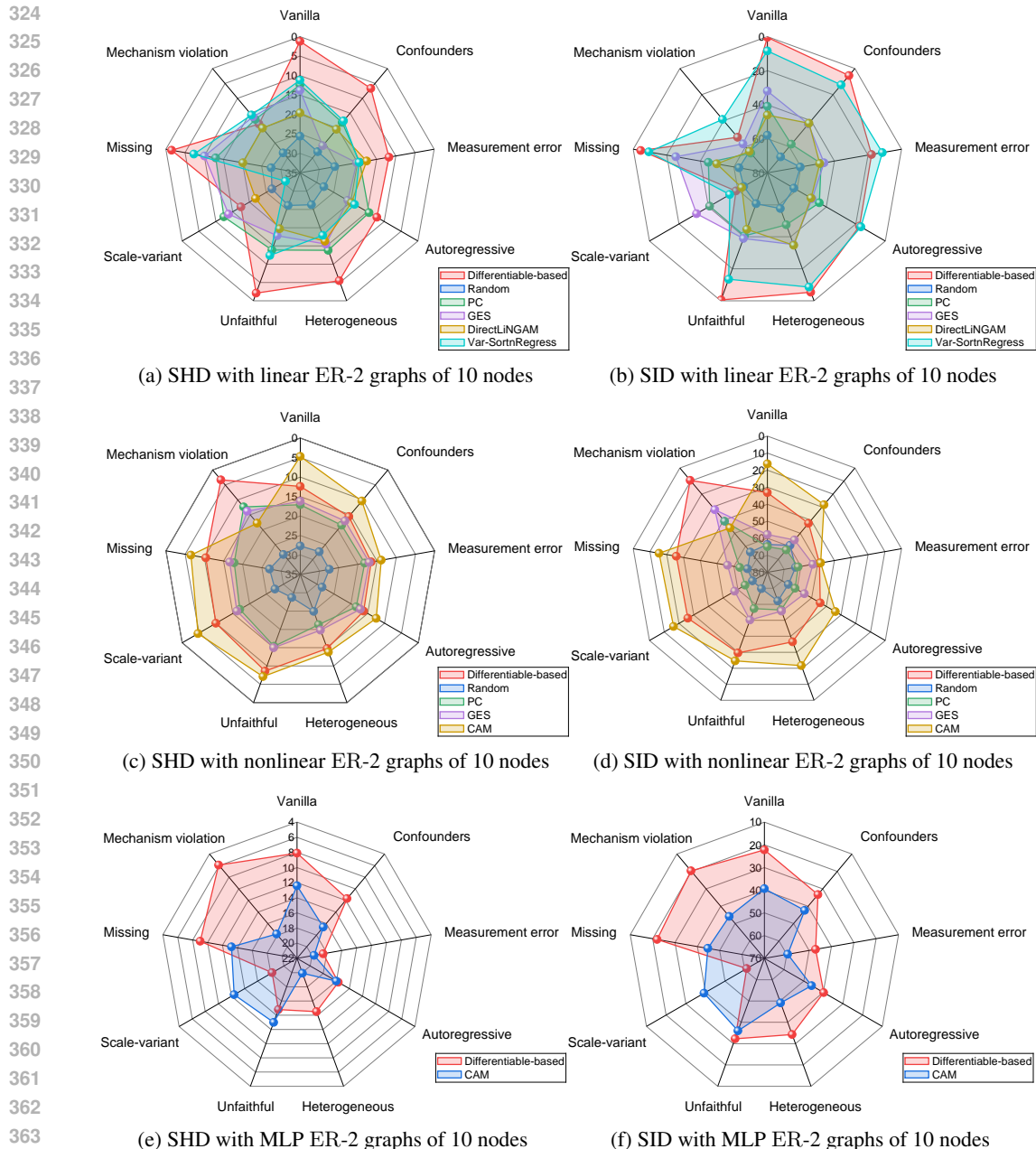


Figure 1: Experimental results under the linear, nonlinear, and mlp settings for both the vanilla scenario and the eight misspecified scenarios. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials on the 10 nodes ER-2 graphs. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 1c and Figure 1d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

**Mechanism violation.** For mechanism violation in the linear setting of Table 2.1 and 2.2, despite PC and GES’s ability to handle nonlinear mechanisms, we are surprised to observe that linear differentiable causal discovery algorithms achieve competitive performance. In the nonlinear setting of Table 3.1 and 3.2, we discover that nonlinear differentiable causal discovery algorithms, such as NOTEARS-MLP (Zheng et al., 2020), Gr $\alpha$ N-DAG (Lachapelle et al., 2019), and DAGMA (Bello et al., 2022), outperform other types of algorithms. From Table 3.1 and 3.2, we also see that CAM does not perform well under mechanism violation, although it excels in other scenarios. We speculate that this is because the Gaussian process mechanism used in the nonlinear vanilla model aligns well

with CAM’s assumptions about SEM. For a more detailed discussion on CAM performance, see Section 4.1.1.

Table 2.1: Linear Setting, for ER-2 graphs of 10 nodes (Part I).

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive	
	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓
Random	25.6±3.1	57.9±9.5	27.9±2.3	67.8±7.8	25.9±3.5	60.4±11.3	27.9±3.2	62.0±8.1
PC	12.4±3.1	40.9±13.4	18.1±4.7	58.1±15.6	19.4±4.1	48.0±13.1	14.5±2.0	44.8±9.5
GES	13.8±7.8	32.0±13.6	25.9±7.7	42.6±14.0	20.2±4.8	46.2±16.7	20.8±5.5	49.7±11.5
DirectLiNGAM	19.6±3.3	46.1±10.6	20.4±5.0	42.0±6.0	17.6±2.4	48.8±12.4	19.7±4.2	50.4±8.4
Var-SortnRegress	11.2±3.5	8.4±8.5	17.6±5.8	12.6±9.9	19.6±2.8	<b>11.4±8.7</b>	18.8±2.4	<b>16.5±10.6</b>
$R^2$ -SortnRegress	20.2±4.8	32.4±14.0	25.7±4.1	37.6±13.0	25.6±6.0	39.2±16.0	25.6±4.9	38.8±19.0
NOTEARS	1.5±1.6	1.8±4.2	8.5±3.9	9.5±8.1	12.5±2.0	19.6±8.6	<b>12.2±3.6</b>	27.5±14.2
GOLEM	1.4±1.4	<b>0.4±1.2</b>	<b>6.7±2.8</b>	14.2±9.8	17.8±2.5	43.1±13.3	16.6±4.0	34.9±16.9
NoCurl	2.0±1.8	5.1±5.8	9.1±4.2	<b>5.4±3.9</b>	<b>11.8±1.8</b>	<b>17.9±8.4</b>	14.8±2.5	<b>17.5±10.8</b>
DAGMA	<b>1.2±1.2</b>	3.3±5.3	8.4±3.9	8.8±7.7	12.6±2.5	18.5±8.6	<b>12.2±3.6</b>	28.4±15.3

Table 2.2: Linear Setting, for ER-2 graphs of 10 nodes (Part II).

10 nodes	Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓
Random	26.3±3.5	57.7±7.6	26.1±3.7	60.7±12.9	26.7±3.0	64.0±8.7	27.5±3.4	63.1±6.0	28.3±3.0	63.6±7.8
PC	13.8±2.6	47.5±10.3	13.9±3.2	40.6±10.4	<b>12.4±3.1</b>	40.9±13.4	13.0±4.7	44.8±16.0	17.1±2.5	64.9±10.0
GES	15.5±6.1	35.2±12.2	17.8±6.4	39.0±15.7	13.8±7.8	<b>32.0±13.6</b>	10.1±5.2	25.4±12.6	16.2±2.2	57.8±10.7
DirectLiNGAM	16.3±3.9	34.7±9.9	19.7±4.3	44.7±13.9	21.8±4.3	62.6±9.3	20.1±4.3	49.8±11.1	20.0±0.0	63.5±7.7
Var-SortnRegress	17.9±3.3	8.6±9.3	12.4±3.1	13.5±8.0	30.7±5.1	54.5±10.3	7.3±3.5	9.3±8.4	<b>15.6±3.3</b>	<b>39.0±6.7</b>
$R^2$ -SortnRegress	26.0±5.4	37.0±14.4	29.8±4.8	51.0±11.3	20.2±4.8	32.4±14.0	20.5±6.7	32.0±8.8	20.3±3.7	66.1±9.7
NOTEARS	<b>5.5±2.7</b>	<b>5.4±5.1</b>	2.7±3.1	3.1±5.1	18.0±1.2	60.5±7.3	2.3±1.7	6.4±8.6	19.0±0.9	58.3±8.0
GOLEM	6.5±4.5	9.8±8.1	<b>2.1±2.2</b>	<b>0.6±1.8</b>	<b>17.5±1.2</b>	64.4±6.8	1.7±1.7	6.2±10.8	<b>18.6±1.6</b>	<b>52.7±4.3</b>
NoCurl	6.6±2.9	5.5±5.7	2.2±2.3	2.0±4.4	27.2±5.1	69.9±7.9	3.1±3.2	4.7±5.8	19.1±1.0	58.9±9.5
DAGMA	<b>5.5±2.3</b>	12.0±8.2	<b>2.1±2.2</b>	<b>0.6±1.8</b>	17.9±1.4	<b>58.7±6.8</b>	<b>1.5±1.4</b>	<b>4.5±7.1</b>	19.0±0.9	58.3±8.0

**Scale-variant model.** In Table 2.1 and 2.2, we observe that the results of linear differentiable causal discovery algorithms, such as NOTEARS (Zheng et al., 2018), GOLEM (Ng et al., 2020), NoCurl (Yu et al., 2021), and DAGMA, significantly decline under scale-variant data, performing worse than PC and GES, which is consistent with the observations of Reisach et al. (2021). For nonlinear differentiable methods, performance under scale-variant data has not been explored in previous research. Table 3.1, 3.2, 4.1 and 4.2 indicate that nonlinear differentiable methods also show performance degradation under scale variation scenarios, and their results almost always lower than CAM. However, unlike the linear scenarios, the result of nonlinear differentiable algorithms is almost always superior to PC and GES.

**Unfaithful model.** In the linear setting of Table 2.1 and 2.2, we see a significant performance drop for Var-SortnRegress (Reisach et al., 2021) and  $R^2$ -SortnRegress (Reisach et al., 2023) under unfaithful distributions. The explanation for this is that for each triplet  $X_i \rightarrow X_j \rightarrow X_k \leftarrow X_i$  in the graph, after the causal direct effect of  $X_i \rightarrow X_k$  cancels out, the variance of node  $X_k$  changes significantly. This reduces the Var-Sortability, further leading to a performance decline in the two SortnRegress algorithms and linear differentiable methods. In the nonlinear setting of Table 3.1 and 3.2, the SHD of various algorithms generally decline to some extent under unfaithful path cancellations. This is consistent with the experimental results of Montagna et al. (2023), which indicate that the cancellation of causal effects in unfaithful nonlinear scenarios makes structural inference of sparse graphs easier.

#### 4.1.1 DISCUSSION ON CAM PERFORMANCE

**Motivations.** The nonlinear vanilla model adopts a Gaussian process that is consistent with the assumptions of the CAM. To provide a fair benchmark for CAM, we consider the nonlinear vanilla model following different functional mechanism. We compare CAM with the representative differentiable causal discovery method: NOTEARS-MLP.

**Simulations.** We simulate ER-2 graphs based on the number of nodes  $d \in \{10, 20, 50\}$ . Following the data generation mechanisms of Zheng et al. (2020), we consider  $f_i$  in nonlinear vanilla model (Section 3.1) is modified to be parameterized by a neural network with one hidden layer of size 100.

**Results.** From Table 4.1 and 4.2, we observe that NOTEARS-MLP achieves better performance under almost all model assumption violations. Considering that the functional mechanisms of data in real-world scenarios are usually unknown, we believe that differentiable causal discovery has



a significant advantage over CAM in all types of assumption violation scenarios except for scale variation.

Table 3.1: Nonlinear Setting, for ER-2 graphs of 10 nodes (Part I).

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive	
	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓
Random	27.7±3.2	63.6±11.2	27.4±2.5	59.3±9.4	27.4±3.6	63.2±7.7	28.5±2.3	65.8±8.6
PC	17.1±2.5	64.9±10.0	18.5±1.9	62.8±9.1	18.2±1.1	61.6±12.2	18.4±1.3	61.2±9.7
GES	16.2±2.2	57.8±10.7	17.1±2.1	55.2±15.5	17.0±1.0	52.8±9.1	17.2±2.0	54.9±7.8
CAM	<b>4.7±1.9</b>	<b>16.3±9.5</b>	<b>10.4±2.8</b>	<b>28.2±5.3</b>	<b>13.9±1.9</b>	<b>48.2±7.4</b>	<b>12.5±3.0</b>	<b>33.9±16.3</b>
NOTEARS-MLP	<b>12.4±2.2</b>	36.3±7.1	17.0±1.7	49.2±8.6	<b>16.5±0.8</b>	<b>48.9±4.8</b>	17.0±3.7	47.7±11.0
GraN-DAG	12.7±2.4	<b>33.2±10.6</b>	<b>15.6±2.1</b>	<b>42.4±8.8</b>	20.0±1.1	63.8±11.3	<b>16.2±2.3</b>	<b>44.2±10.0</b>
DAGMA	13.5±2.0	40.7±8.1	18.6±2.0	62.0±13.3	17.3±1.3	54.9±8.3	19.0±2.0	56.6±10.5

Table 3.2: Nonlinear Setting, for ER-2 graphs of 10 nodes (Part II).

10 nodes	Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓
Random	24.9±3.2	62.1±10.1	28.7±2.1	69.8±7.9	27.5±2.5	69.9±5.9	27.0±4.2	68.1±8.0	28.3±3.4	64.5±10.7
PC	21.3±3.2	56.5±10.4	15.4±1.4	57.3±9.5	17.1±2.5	64.9±10.0	17.8±3.1	63.5±10.5	12.4±3.1	40.9±13.4
GES	19.8±2.5	55.7±8.4	15.0±4.5	50.2±13.1	16.2±2.2	57.8±10.7	16.6±2.6	56.2±9.9	13.8±7.8	32.0±13.6
CAM	<b>13.8±2.9</b>	<b>21.7±12.8</b>	<b>7.1±3.2</b>	<b>24.7±13.0</b>	<b>4.7±1.9</b>	<b>16.3±9.5</b>	<b>6.5±2.1</b>	<b>15.4±7.3</b>	17.8±4.4	45.9±17.0
NOTEARS-MLP	16.4±3.7	42.6±10.6	11.4±2.1	43.8±9.2	16.1±2.5	48.3±9.7	12.3±2.1	33.6±7.0	5.9±2.5	19.7±8.7
GraN-DAG	<b>14.8±2.1</b>	40.5±10.8	11.2±3.1	37.7±10.0	<b>10.0±3.7</b>	<b>26.1±12.0</b>	<b>10.4±3.4</b>	<b>25.7±8.1</b>	16.3±2.1	54.7±5.1
DAGMA	16.4±4.3	<b>36.6±15.7</b>	<b>8.6±3.1</b>	<b>29.6±14.8</b>	15.7±2.7	53.3±12.9	13.7±2.7	41.4±7.8	<b>3.3±3.1</b>	<b>9.5±11.6</b>

#### 4.1.2 THEORY ON DIFFERENTIABLE CAUSAL DISCOVERY IN MISSPECIFIED SCENARIOS

We analyze the performance of linear differentiable causal discovery under measurement error, unfaithful and missing scenarios by introducing the theories (Theorem 7 and Theorem 9) from Loh & Bühlmann (2014). Theorem 7 in Loh & Bühlmann (2014) states that for a linear model with equal noise variance, minimizing the least squares score will return the true DAG in the large sample limit. For a linear model with non-equal noise variances, we define the noise ratio

$$r = \frac{\max(\sigma_1^2, \dots, \sigma_d^2)}{\min(\sigma_1^2, \dots, \sigma_d^2)}. \quad (6)$$

Theorem 9 in Loh & Bühlmann (2014) states that if  $r < 1 + \frac{\xi}{d}$ , where  $\xi > 0$  is the gap between the score of the true DAG and the next best DAG, minimizing the least squares score will return the true DAG in the large sample limit.

**Measurement error.** In measurement error model considered by Montagna et al. (2023), the observed variables are:

$$\tilde{X}_i = X_i + \epsilon_i, \forall i = 1, \dots, d, \quad (7)$$

where  $X_i = f_i(X_{pa}(X_i)) + U_i$ ,  $f_i$  is a linear mechanism,  $U_i \sim N(0, 1)$ ,  $\epsilon_i \sim N(0, \delta * \text{Var}(X_i))$  with  $\delta \in \{0.2, 0.4, 0.6, 0.8\}$ . In the vanilla model,  $r = 1$ . However, in the measurement error model, the noise ratio becomes

$$\tilde{r} = \frac{\max(1 + \delta * \text{Var}(X_i), \dots, 1 + \delta * \text{Var}(X_d))}{\min(1 + \delta * \text{Var}(X_i), \dots, 1 + \delta * \text{Var}(X_d))}. \quad (8)$$

Due to the increasing trend of marginal variances of nodes along the causal direction (Reisach et al., 2021), we infer that  $\tilde{r} > r = 1$ . In this scenario, there is no guarantee that  $\tilde{r} < 1 + \frac{\xi}{d}$  and linear differentiable causal discovery based on least squares cannot guarantee obtaining the true DAG, which can explain their performance decline in Table 2.1 and 2.2.

**Unfaithful model.** In unfaithful model considered by Montagna et al. (2023), for each triplet  $X_i \rightarrow X_j \rightarrow X_k \leftarrow X_i$  in the graph, the causal mechanisms are adjusted such that the direct effect of  $X_i$  on  $X_k$  cancels out. To illustrate the change in noise ratio after path cancellation, we consider a DAG  $\mathcal{G}$  with variable set  $V(\mathcal{G}) = \{X_1, X_2, X_3\}$  and edge set  $E(\mathcal{G}) = \{X_1 \rightarrow X_2, X_2 \rightarrow X_3, X_1 \rightarrow X_3\}$ . The structural equations is defined as:

$$\begin{aligned} X_1 &= U_1, \\ X_2 &= f_1(X_1) + U_2, \\ X_3 &= f_1(X_1) - X_2 + U_3, \end{aligned} \quad (9)$$

where  $f_1$  is a linear mechanism. After the direct causal effect of  $X_1 \rightarrow X_3$  cancels out, the noise term of  $X_3$  is  $U_3 - U_2$  with the distribution of  $N(0, 2)$ . Similarly, in the unfaithful datasets with nodes  $d \in \{10, 20, 50\}$  considered by our experiments, for each triplet  $X_i \rightarrow X_j \rightarrow X_k \leftarrow X_i$  in the graph, once unfaithful path cancellation occurs, the noise term of  $X_k$  is  $U_k - U_j$  with the distribution of  $N(0, 2)$ . In this case, the noise ratio becomes  $r' = 2 > r = 1$  (vanilla model). Due to the increasing of the noise ratio, there is no guarantee that  $r' < 1 + \frac{\xi'_d}{d}$  and linear differentiable causal discovery based on least squares cannot guarantee obtaining the true DAG, which can also explain their performance decline in Table 2.1 and 2.2.

**Missing model.** Under the MCAR case, we deleted the rows with missing values and regenerated the data under the i.i.d. assumption to ensure the unchanged sample size. The noise ratio  $r = 1$  remains constant before and after data imputation. This explains the superior performance of differentiable causal discovery methods observed in Table 2.1 and 2.2.

Table 4.1: MLP Setting, for ER-2 graphs of 10 nodes (Part I).

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive	
	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓
CAM	12.4±3.6	39.3±16.5	16.6±4.2	42.4±17.7	19.7±4.7	59.6±12.0	16.0±3.4	46.0±16.6
NOTEARS-MLP	<b>8.1±2.7</b>	<b>22.2±10.6</b>	<b>11.7±5.5</b>	<b>33.3±17.0</b>	<b>18.5±3.7</b>	<b>47.1±11.9</b>	<b>15.7±4.3</b>	<b>39.8±9.6</b>

Table 4.2: MLP Setting, for ER-2 graphs of 10 nodes (Part II).

10 nodes	Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓	SHD↓	SID↓
CAM	19.9±3.6	49.1±7.2	<b>13.0±3.2</b>	36.0±14.0	<b>12.4±3.6</b>	<b>39.3±16.5</b>	13.2±3.6	44.7±16.7	17.8±4.4	45.9±17.0
NOTEARS-MLP	<b>14.5±2.2</b>	<b>34.3±10.4</b>	14.8±3.9	<b>32.3±12.8</b>	18.2±2.6	61.0±11.0	<b>9.0±3.3</b>	<b>22.0±10.0</b>	<b>5.9±2.5</b>	<b>19.7±8.7</b>

## 4.2 SUMMARY AND IMPLICATIONS FOR PRACTICE

In Appendix H, we summarize the results of the most competitive methods under misspecified scenarios. Differentiable causal discovery methods demonstrate optimal or competitive performance in commonly used scenarios other than scale variation. Notably, the recent work by Deng et al. (2024) shows that for linear differentiable methods, scale invariance can be achieved by appropriately choosing the loss function. This further reinforces our conclusion regarding the robustness of differentiable methods. In our benchmarks, the results in Table 11, Table 21 and Table 25 indicate that the performance of nonlinear differentiable methods under scale variation remains challenging and warrants further investigation. In practice, the misspecified scenarios are inevitably encountered, making the robustness of algorithms critically important. Based on the summarized results on eight misspecified synthetic datasets (see Appendix H), runtime results of benchmark methods (see Appendix D), real-world (see Appendix I) and semi-synthetic data (see Appendix L) results, we observe that differentiable causal discovery methods have the potential to achieve optimal or competitive performance on real-world data with an almost negligible time cost. The fast and robust characteristics of differentiable methods enable them to better address the challenges of applying causal discovery algorithms to real-world data, demonstrating their practical implementation potential.

## 5 CONCLUSION

This work assesses the efficacy of twelve preeminent causal discovery methods across eight scenarios involving violations of model assumptions. These methods encompass approaches grounded in independence constraints, scoring criteria, functional causal models, and differentiable causal discovery. Our experimental results show that differentiable causal discovery methods exhibit remarkable resilience in commonly used scenarios of model assumption violations, except for scale variation. It is not our intention to assert that differentiable causal discovery will achieve optimal performance across all circumstances, rather, we aim to underscore its substantial potential within the benchmarks we have evaluated, thereby emphasizing the necessity for further exploration in this direction. In future work, causal discovery methods for more semi-synthetic data and real-world scenarios will be explored. Finally, our study confines itself to non-temporal causal discovery algorithms. Equally crucial is the conduct of benchmark assessments for causal discovery in time series and event sequences under model assumption violations.

## REFERENCES

- 540  
541  
542 Hirotoyu Akaike. Information theory and an extension of the maximum likelihood principle. In  
543 *Selected papers of hirotugu akaike*, pp. 199–213. Springer, 1998.
- 544 Kevin Bello, Bryon Aragam, and Pradeep Ravikumar. Dagma: Learning dags via m-matrices and  
545 a log-determinant acyclicity characterization. In *Advances in Neural Information Processing*  
546 *Systems*, volume 35, pp. 8226–8239, 2022.
- 547  
548 Ulrik Brandes, Marco Gaertler, and Dorothea Wagner. Experiments on graph clustering algorithms.  
549 In *European symposium on algorithms*, pp. 568–579. Springer, 2003.
- 550 Philippe Brouillard, Sébastien Lachapelle, Alexandre Lacoste, Simon Lacoste-Julien, and Alexan-  
551 dre Drouin. Differentiable causal discovery from interventional data. In *Advances in Neural*  
552 *Information Processing Systems*, volume 33, pp. 21865–21877, 2020.
- 553  
554 Peter Bühlmann, Jonas Peters, and Jan Ernest. Cam: Causal additive models, high-dimensional  
555 order search and penalized regression. *The Annals of Statistics*, 42(6):2526–2556, 2014. doi:  
556 10.1214/14-AOS1260.
- 557 Ruichu Cai, Feng Xie, Clark Glymour, Zhifeng Hao, and Kun Zhang. Triad constraints for learning  
558 causal structure of latent variables. In *Advances in Neural Information Processing Systems*,  
559 volume 32, 2019.
- 560  
561 Ruichu Cai, Jincheng Ye, Jie Qiao, Huiyuan Fu, and Zhifeng Hao. Fom: Fourth-order moment based  
562 causal direction identification on the heteroscedastic data. *Neural Networks*, 124:193–201, 2020.
- 563  
564 Ruichu Cai, Zhiyi Huang, Wei Chen, Zhifeng Hao, and Kun Zhang. Causal discovery with latent  
565 confounders based on higher-order cumulants. In *International Conference on Machine Learning*,  
566 pp. 3380–3407. PMLR, 2023.
- 567  
568 David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of Machine*  
*Learning Research*, 3(Nov):507–554, 2002.
- 569  
570 Haoyue Dai, Peter Spirtes, and Kun Zhang. Independence testing-based approach to causal discovery  
571 under measurement error and linear non-gaussian models. *Advances in Neural Information*  
572 *Processing Systems*, 35:27524–27536, 2022.
- 573  
574 Chang Deng, Kevin Bello, Pradeep Ravikumar, and Bryon Aragam. Likelihood-based differentiable  
structure learning. *arXiv preprint arXiv:2410.06163*, 2024.
- 575  
576 José A Gámez, Juan L Mateo, and José M Puerta. Learning bayesian networks by hill climbing: effi-  
577 cient methods based on progressive restriction of the neighborhood. *Data Mining and Knowledge*  
578 *Discovery*, 22:106–148, 2011.
- 579  
580 Erdun Gao, Junjia Chen, Li Shen, Tongliang Liu, Mingming Gong, and Howard Bondell. Feddag:  
Federated dag structure learning. *arXiv preprint arXiv:2112.03555*, 2021.
- 581  
582 Erdun Gao, Ignavier Ng, Mingming Gong, Li Shen, Wei Huang, Tongliang Liu, Kun Zhang, and  
583 Howard Bondell. Missdag: Causal discovery in the presence of missing data with continuous  
584 additive noise models. In *Advances in Neural Information Processing Systems*, volume 35, pp.  
585 5024–5038, 2022.
- 586  
587 AmirEmad Ghassami, Saber Salehkaleybar, Negar Kiyavash, and Kun Zhang. Learning causal  
588 structures using regression invariance. In *Advances in Neural Information Processing Systems*,  
volume 30, 2017.
- 589  
590 AmirEmad Ghassami, Negar Kiyavash, Biwei Huang, and Kun Zhang. Multi-domain causal structure  
591 learning in linear systems. In *Advances in Neural Information Processing Systems*, volume 31,  
592 2018.
- 593  
Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on  
graphical models. *Frontiers in Genetics*, 10:524, 2019.

- 594 Alain Hauser and Peter Bühlmann. Characterization and greedy learning of interventional markov  
595 equivalence classes of directed acyclic graphs. *Journal of Machine Learning Research*, 13(1):  
596 2409–2464, 2012.
- 597 Yue He, Peng Cui, Zheyang Shen, Renzhe Xu, Furui Liu, and Yong Jiang. Daring: Differentiable  
598 causal discovery with residual independence. In *Proceedings of the 27th ACM SIGKDD Conference*  
599 *on Knowledge Discovery and Data Mining*, pp. 596–605, 2021.
- 600 David Heckerman, Dan Geiger, and David M Chickering. Learning bayesian networks: The combi-  
601 nation of knowledge and statistical data. *Machine learning*, 20:197–243, 1995.
- 602 Christina Heinze-Deml, Marloes H Maathuis, and Nicolai Meinshausen. Causal structure learning.  
603 *Annual Review of Statistics and Its Application*, 5(1):371–391, 2018.
- 604 Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear  
605 causal discovery with additive noise models. In *Advances in Neural Information Processing*  
606 *Systems*, volume 21, 2008.
- 607 Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour,  
608 and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of*  
609 *Machine Learning Research*, 21(89):1–53, 2020.
- 610 Diviyani Kalainathan and Olivier Goudet. Causal discovery toolbox: Uncover causal relationships in  
611 python. *arXiv preprint arXiv:1903.02278*, 2019.
- 612 Diviyani Kalainathan, Olivier Goudet, Isabelle Guyon, David Lopez-Paz, and Michèle Sebag. Struc-  
613 tural agnostic modeling: Adversarial learning of causal graphs. *Journal of Machine Learning*  
614 *Research*, 23(219):1–62, 2022.
- 615 Mehmet Kayaalp and Gregory F Cooper. A bayesian network scoring metric that is based on globally  
616 uniform parameter priors. *arXiv preprint arXiv:1301.0576*, 2012.
- 617 Lingjing Kong, Biwei Huang, Feng Xie, Eric Xing, Yuejie Chi, and Kun Zhang. Identification of  
618 nonlinear latent hierarchical models. In *Advances in Neural Information Processing Systems*,  
619 volume 36, pp. 2010–2032, 2023.
- 620 Trent Kyono, Yao Zhang, and Mihaela van der Schaar. Castle: Regularization via auxiliary causal  
621 graph discovery. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1501–  
622 1512, 2020.
- 623 Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based  
624 neural dag learning. *arXiv preprint arXiv:1906.02226*, 2019.
- 625 Fangfu Liu, Wenchang Ma, An Zhang, Xiang Wang, Yueqi Duan, and Tat-Seng Chua. Discovering  
626 dynamic causal space for dag structure learning. In *Proceedings of the 29th ACM SIGKDD*  
627 *Conference on Knowledge Discovery and Data Mining*, pp. 1429–1440, 2023.
- 628 Po-Ling Loh and Peter Bühlmann. High-dimensional learning of linear causal networks via inverse  
629 covariance estimation. *Journal of Machine Learning Research*, 15(1):3065–3105, 2014.
- 630 Francesco Montagna, Atalanti Mastakouri, Elias Eulig, Nicoletta Noceti, Lorenzo Rosasco, Dominik  
631 Janzing, Bryon Aragam, and Francesco Locatello. Assumption violations in causal discovery  
632 and the robustness of score matching. In *Advances in Neural Information Processing Systems*,  
633 volume 36, 2023.
- 634 Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. Distinguishing  
635 cause from effect using observational data: methods and benchmarks. *Journal of Machine*  
636 *Learning Research*, 17(32):1–102, 2016.
- 637 Joris M Mooij, Sara Magliacane, and Tom Claassen. Joint causal inference from multiple contexts.  
638 *Journal of Machine Learning Research*, 21(99):1–108, 2020.
- 639 Andrew A Neath and Joseph E Cavanaugh. The bayesian information criterion: background,  
640 derivation, and applications. *Wiley Interdisciplinary Reviews: Computational Statistics*, 4(2):  
641 199–203, 2012.

- 648 AS Nemirovsky. Optimization ii. numerical methods for nonlinear continuous optimization. 1999.  
649
- 650 Ignavier Ng, Shengyu Zhu, Zhitang Chen, and Zhuangyan Fang. A graph autoencoder approach to  
651 causal structure learning. *arXiv preprint arXiv:1911.07420*, 2019.
- 652 Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints  
653 for learning linear dags. In *Advances in Neural Information Processing Systems*, volume 33, pp.  
654 17943–17954, 2020.
- 655
- 656 Ignavier Ng, Biwei Huang, and Kun Zhang. Structure learning with continuous optimization: A  
657 sober look and beyond. In *Causal Learning and Reasoning*, pp. 71–105. PMLR, 2024.
- 658 Roxana Pamfil, Nisara Sriwattanaworachai, Shaan Desai, Philip Pilgerstorfer, Konstantinos Geor-  
659 gatzis, Paul Beaumont, and Bryon Aragam. Dynotears: Structure learning from time-series data.  
660 In *International Conference on Artificial Intelligence and Statistics*, pp. 1595–1605. PMLR, 2020.  
661
- 662 Judea Pearl. *Causality*. Cambridge university press, 2009.
- 663 Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic books,  
664 2018.
- 665
- 666 Jonas Peters and Peter Bühlmann. Identifiability of gaussian structural equation models with equal  
667 error variances. *Biometrika*, 101(1):219–228, 2014.
- 668 Jonas Peters and Peter Bühlmann. Structural intervention distance for evaluating causal graphs.  
669 *Neural Computation*, 27(3):771–799, 2015.  
670
- 671 Jonas Peters, Joris M. Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery with  
672 continuous additive noise models. *Journal of Machine Learning Research*, 15(58):2009–2053,  
673 2014. URL <http://jmlr.org/papers/v15/peters14a.html>.
- 674 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations  
675 and learning algorithms*. The MIT Press, 2017.  
676
- 677 Alexander Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag! causal  
678 discovery benchmarks may be easy to game. In *Advances in Neural Information Processing  
679 Systems*, volume 34, pp. 27772–27784, 2021.
- 680 Alexander Reisach, Myriam Tami, Christof Seiler, Antoine Chambaz, and Sebastian Weichwald. A  
681 scale-invariant sorting criterion to find a causal order in additive noise models. In *Advances in  
682 Neural Information Processing Systems*, volume 36, 2023.
- 683
- 684 Robert W Robinson. Counting labeled acyclic digraphs. *New directions in the theory of graphs*, pp.  
685 239–273, 1973.
- 686 Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal protein-  
687 signaling networks derived from multiparameter single-cell data. *Science*, 308(5721):523–529,  
688 2005.
- 689
- 690 Saber Salehkaleybar, AmirEmad Ghassami, Negar Kiyavash, and Kun Zhang. Learning linear non-  
691 gaussian causal models in the presence of latent variables. *Journal of Machine Learning Research*,  
692 21(39):1–24, 2020.
- 693 Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. A linear  
694 non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10),  
695 2006.
- 696
- 697 Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara,  
698 Takashi Washio, Patrik O Hoyer, Kenneth Bollen, and Patrik Hoyer. Directlingam: A direct  
699 method for learning a linear non-gaussian structural equation model. *Journal of Machine Learning  
700 Research*, 12(Apr):1225–1248, 2011.
- 701 Karamjit Singh, Garima Gupta, Vartika Tewari, and Gautam Shroff. Comparative benchmarking of  
causal discovery techniques. *arXiv preprint arXiv:1708.06246*, 2017.

- 702 Peter Spirtes. Introduction to causal inference. *Journal of Machine Learning Research*, 11(5), 2010.  
703
- 704 Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social*  
705 *science computer review*, 9(1):62–72, 1991.
- 706 Peter Spirtes, Christopher Meek, and Thomas Richardson. Causal inference in the presence of latent  
707 variables and selection bias. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial*  
708 *Intelligence*, pp. 499–506, 1995.
- 709 Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press,  
710 2001.
- 711
- 712 Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing bayesian  
713 network structure learning algorithm. *Machine learning*, 65:31–78, 2006.
- 714 Ruibo Tu, Cheng Zhang, Paul Ackermann, Karthika Mohan, Hedvig Kjellström, and Kun Zhang.  
715 Causal discovery in the presence of missing data. In *The 22nd International Conference on*  
716 *Artificial Intelligence and Statistics*, pp. 1762–1770. PMLR, 2019a.
- 717
- 718 Ruibo Tu, Kun Zhang, Bo Bertilson, Hedvig Kjellstrom, and Cheng Zhang. Neuropathic pain  
719 diagnosis simulator for causal discovery algorithm evaluation. In *Advances in Neural Information*  
720 *Processing Systems*, volume 32, 2019b.
- 721 Matthew J Vowels, Necati Cihan Camgoz, and Richard Bowden. D’ya like dags? a survey on  
722 structure learning and causal discovery. *ACM Computing Surveys*, 55(4):1–36, 2022.
- 723
- 724 Yu Wang, An Zhang, Xiang Wang, Yancheng Yuan, Xiangnan He, and Tat-Seng Chua. Differentiable  
725 invariant causal discovery. *arXiv preprint arXiv:2205.15638*, 2022.
- 726
- 727 Dennis Wei, Tian Gao, and Yue Yu. Dags with no fears: A closer look at continuous optimization for  
728 learning bayesian networks. *Advances in Neural Information Processing Systems*, 33:3895–3906,  
729 2020.
- 730 Feng Xie, Ruichu Cai, Biwei Huang, Clark Glymour, Zhifeng Hao, and Kun Zhang. Generalized  
731 independent noise condition for estimating latent variable causal graphs. In *Advances in Neural*  
732 *Information Processing Systems*, volume 33, pp. 14891–14902, 2020.
- 733 Shuai Yang, Kui Yu, Fuyuan Cao, Lin Liu, Hao Wang, and Jiuyong Li. Learning causal representations  
734 for robust domain adaptation. *IEEE Transactions on Knowledge and Data Engineering*, 35(3):  
735 2750–2764, 2021.
- 736 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks.  
737 In *International Conference on Machine Learning*, pp. 7154–7163. PMLR, 2019.
- 738
- 739 Yue Yu, Tian Gao, Naiyu Yin, and Qiang Ji. Dags with no curl: An efficient dag structure learning  
740 approach. In *International Conference on Machine Learning*, pp. 12156–12166. PMLR, 2021.
- 741
- 742 An Zhang, Fangfu Liu, Wenchang Ma, Zhibo Cai, Xiang Wang, and Tat-Seng Chua. Boosting  
743 differentiable causal discovery via adaptive sample reweighting. *arXiv preprint arXiv:2303.03187*,  
744 2023.
- 745 Keli Zhang, Shengyu Zhu, Marcus Kalander, Ignavier Ng, Junjian Ye, Zhitang Chen, and Lujia Pan.  
746 gcastle: A python toolbox for causal discovery. *arXiv preprint arXiv:2111.15155*, 2021.
- 747 Kun Zhang and Aapo Hyvarinen. On the identifiability of the post-nonlinear causal model. *arXiv*  
748 *preprint arXiv:1205.2599*, 2012.
- 749
- 750 Kun Zhang, Mingming Gong, Joseph Ramsey, K. Batmanghelich, Peter Spirtes, and Clark Gly-  
751 mour. Causal discovery with linear non-gaussian models under measurement error: Structural  
752 identifiability results. In *Conference on Uncertainty in Artificial Intelligence*, 2018. URL  
753 <https://api.semanticscholar.org/CorpusID:54058643>.
- 754 Muhan Zhang, Shali Jiang, Zhicheng Cui, Roman Garnett, and Yixin Chen. D-vae: A variational  
755 autoencoder for directed acyclic graphs. In *Advances in Neural Information Processing Systems*,  
volume 32, 2019.

756 Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous  
757 optimization for structure learning. In *Advances in Neural Information Processing Systems*,  
758 volume 31, 2018.

759 Xun Zheng, Chen Dan, Bryon Aragam, Pradeep Ravikumar, and Eric Xing. Learning sparse  
760 nonparametric dags. In *International Conference on Artificial Intelligence and Statistics*, pp.  
761 3414–3425. PMLR, 2020.

762 Yujia Zheng, Biwei Huang, Wei Chen, Joseph Ramsey, Mingming Gong, Ruichu Cai, Shohei Shimizu,  
763 Peter Spirtes, and Kun Zhang. Causal-learn: Causal discovery in python. *Journal of Machine*  
764 *Learning Research*, 25(60):1–8, 2024.

765 Rong Zhu, Andreas Pfadler, Ziniu Wu, Yuxing Han, Xiaoke Yang, Feng Ye, Zhenping Qian, Jingren  
766 Zhou, and Bin Cui. Efficient and scalable structure learning for bayesian networks: Algorithms  
767 and applications. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pp.  
768 2613–2624. IEEE, 2021.

769 Shengyu Zhu, Ignavier Ng, and Zhitang Chen. Causal discovery with reinforcement learning. *arXiv*  
770 *preprint arXiv:1906.04477*, 2019.

771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

810	CONTENTS	
811		
812	<b>1 Introduction</b>	<b>1</b>
813		
814	<b>2 Background</b>	<b>3</b>
815		
816	2.1 Task formulation . . . . .	3
817		
818	2.2 Structure identifiability . . . . .	3
819		
820	2.3 Differentiable score-based causal discovery . . . . .	3
821		
822	<b>3 Experimental design</b>	<b>4</b>
823		
824	3.1 Synthetic datasets . . . . .	4
825	3.1.1 Model assumption violation scenarios . . . . .	5
826	3.1.2 Data generation . . . . .	5
827		
828	3.2 Methods . . . . .	5
829		
830	3.3 Hyperparameter settings . . . . .	6
831		
832	3.4 Evaluation metrics . . . . .	6
833		
834	<b>4 Critical experimental results and insights</b>	<b>6</b>
835		
836	4.1 Current methods' performance in misspecified scenarios . . . . .	6
837	4.1.1 Discussion on CAM performance . . . . .	8
838	4.1.2 Theory on differentiable causal discovery in misspecified scenarios . . . . .	9
839		
840	4.2 Summary and implications for practice . . . . .	10
841		
842	<b>5 Conclusion</b>	<b>10</b>
843		
844	<b>A Causal Assumptions</b>	<b>18</b>
845		
846	A.1 Causal Markov Property . . . . .	18
847		
848	A.2 Faithfulness . . . . .	18
849		
850	A.3 Causal Sufficiency . . . . .	18
851		
852	A.4 Independent and identically distributed . . . . .	18
853		
854	A.5 Equal Noise Variances . . . . .	18
855		
856	<b>B Benchmark methods</b>	<b>19</b>
857		
858	B.1 PC . . . . .	19
859		
860	B.2 GES . . . . .	19
861		
862	B.3 DirectLiNGAM . . . . .	19
863		
	B.4 CAM . . . . .	19
	B.5 SortnRegress . . . . .	19
	B.6 NOTEARS . . . . .	20
	B.7 GOLEM . . . . .	20
	B.8 NOTEARS-MLP . . . . .	21



864	B.9 GraN-DAG . . . . .	22
865	B.10 NOCURL . . . . .	22
866	B.11 DAGMA . . . . .	23
867		
868		
869	<b>C Related work</b>	<b>23</b>
870		
871	<b>D Table results for runtime of the benchmark methods</b>	<b>24</b>
872		
873	<b>E Table results across nodes, graph types, and graph densities</b>	<b>25</b>
874		
875	<b>F Table results for combined misspecified scenarios</b>	<b>28</b>
876		
877	<b>G Table results for non-Gaussian noise</b>	<b>28</b>
878		
879	<b>H Summary of the most competitive methods</b>	<b>29</b>
880		
881	<b>I Table Results on real-world data</b>	<b>30</b>
882		
883	<b>J Figure results across nodes, graph types, and graph densities</b>	<b>31</b>
884		
885	<b>K Table results for extreme measurement error</b>	<b>38</b>
886		
887	<b>L Table results on semi-synthetic data</b>	<b>39</b>
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		

## 918 A CAUSAL ASSUMPTIONS

919 In this section, we introduce the assumptions frequently used in the causal discovery literature.

### 922 A.1 CAUSAL MARKOV PROPERTY

923 The joint probability distribution  $P(X)$  satisfies the global Markov property (Peters et al., 2017) with respect to the DAG  $\mathcal{G}$  if

$$924 X_A \perp\!\!\!\perp_{\mathcal{G}} X_B \mid X_C \Rightarrow X_A \perp\!\!\!\perp_{P(X)} X_B \mid X_C, \quad (10)$$

925 where  $X_A$ ,  $X_B$  and  $X_C$  are the disjoint subsets of  $X = (X_1, \dots, X_d) \in \mathbb{R}^d$ ,  $\perp\!\!\!\perp_{\mathcal{G}}$  denotes *d-separation* in the causal graph  $\mathcal{G}$ , and  $\perp\!\!\!\perp_{P(X)}$  represents independence in the joint probability distribution  $P(X)$ .

931 In the causal graph  $\mathcal{G}$  of SCM, each variable is independent of its non-descendant nodes when its parents are known, which is referred to as the local Markov property (Peters et al., 2017). Causal Markov property also implies that the joint probability distribution  $P(X)$  can be factorized in the following form:

$$932 P(X) = \prod_i^d P(X_i \mid pa(X_i)). \quad (11)$$

### 938 A.2 FAITHFULNESS

939 The joint probability distribution  $P(X)$  is faithful (Peters et al., 2017) to the DAG  $\mathcal{G}$  if

$$940 X_A \perp\!\!\!\perp_{P(X)} X_B \mid X_C \Rightarrow X_A \perp\!\!\!\perp_{\mathcal{G}} X_B \mid X_C. \quad (12)$$

941 Faithfulness assumption implies that the conditional independence in the  $P(X)$  can be used to infer the graph structure. Constraint-based and traditional score-based causal discovery are typically founded on the faithfulness assumption. In a causal graph, the cancellation of effects along multiple causal paths can lead to a violation of faithfulness assumption.

### 947 A.3 CAUSAL SUFFICIENCY

948 Causal sufficiency assumption is also referred to as no latent confounder. A set of variables  $X$  is said to satisfy the causal sufficiency if there is no unobserved common cause variable  $C$  that influences more than one variable in  $X$  (Spirtes, 2010). This assumption is also frequently considered in causal discovery literature. However, since we cannot always observe all variables in the real world, causal sufficiency assumption is inevitably violated.

### 954 A.4 INDEPENDENT AND IDENTICALLY DISTRIBUTED

955 Non-temporal causal discovery algorithms typically also require the i.i.d. assumption. In the main text, heterogeneous multi-domain data and autoregressive scenarios are two special cases where the i.i.d. assumption is violated. Heterogeneous multi-domain data is closely related to the non-stationary time series data considered in the literature on temporal causal discovery (Huang et al., 2020). Below, we introduce the connection between them. We consider the distribution of  $X_i$  changing with domain or time index, where the mechanism for the  $t$ -th data point is as follows:

$$956 X_{i,t} = f_{i,t}(pa(X_{i,t}), \epsilon_{i,t}), \quad (13)$$

957 where  $\epsilon_{i,t}$  is the noise term of  $X_{i,t}$ . In heterogeneous multi-domain data,  $t$  represents the domain index, whereas in non-stationary time series data,  $t$  denotes the time index.

### 966 A.5 EQUAL NOISE VARIANCES

967 Ng et al. (2024) observe that the performance of linear differentiable causal discovery methods significantly declines in data with non-equal noise variances. They hypothesize that this may be due to the optimization problem becoming severely non-convex under non-equal noise variances, leading to local optimal solutions. Although differentiable methods do not explicitly assume equal noise variance, the performance decline suggests treating equal noise variance as a causal assumption.

## 972 B BENCHMARK METHODS

### 973 B.1 PC

974 PC (Spirtes & Glymour, 1991) algorithm is a representative constraint-based causal discovery method.  
 975 In the first step, under the faithfulness assumption, the global causal skeleton is determined based  
 976 on conditional independence tests. In the second step, some edge directions in the skeleton are  
 977 determined by identifying collider structures. Finally, the remaining edge directions are determined  
 978 using orientation rules, resulting in the MEC. We use the implementation of the PC algorithm in  
 979 causal-learn (Zheng et al., 2024) python package, available at [https://github.com/py-why/](https://github.com/py-why/causal-learn)  
 980 [causal-learn](https://github.com/py-why/causal-learn).  
 981  
 982

### 983 B.2 GES

984 GES (Chickering, 2002) algorithm is a classical score-based causal discovery method. GES mainly  
 985 includes two stages. In the first stage, starting from an empty graph, edges are added through  
 986 greedy equivalence search, and the structures in the equivalence class of the new graph are scored.  
 987 The graph with the highest score is selected, and the edge-adding process is repeated until the  
 988 score reaches a local maximum. In the second stage, starting from the graph obtained in the first  
 989 stage, edges are removed through greedy equivalence search, and the structures in the equivalence  
 990 class of the new graph are scored. The graph with the highest score is selected, and the edge-  
 991 removal process is repeated until the score reaches a local maximum. We use the implementation  
 992 of the GES algorithm in causal-learn (Zheng et al., 2024) python package, available at [https://github.com/py-why/](https://github.com/py-why/causal-learn)  
 993 [causal-learn](https://github.com/py-why/causal-learn).  
 994  
 995

### 996 B.3 DIRECTLINGAM

997 DirectLiNGAM (Shimizu et al., 2011) is a classical linear method based on the functional  
 998 causal model. To address the issues of slow convergence and large errors in the ICA-based  
 999 LiNGAM (Shimizu et al., 2006) algorithm, Shimizu et al. (2011) proposed DirectLiNGAM based  
 1000 on the principle of residual independence. Although the solving speed became slower, the accuracy  
 1001 and convergence improved. We use the implementation of DirectLiNGAM algorithm in gCas-  
 1002 tle (Zhang et al., 2021) python package, available at [https://github.com/huawei-noah/](https://github.com/huawei-noah/trustworthyAI/tree/master/gcastle)  
 1003 [trustworthyAI/tree/master/gcastle](https://github.com/huawei-noah/trustworthyAI/tree/master/gcastle).  
 1004

### 1005 B.4 CAM

1006 CAM (Bühlmann et al., 2014) is a method used for high-dimensional additive structural equation  
 1007 models. CAM separates the search for the order of variables from the selection of edges. It  
 1008 performs the variable order search through nonregularized maximum likelihood estimation and uses  
 1009 sparse regression techniques for edge selection. We use the implementation of CAM algorithm  
 1010 in Causal Discovery Toolbox (Kalainathan & Goudet, 2019) python package, available at [https://github.com/](https://github.com/FenTechSolutions/CausalDiscoveryToolbox)  
 1011 <https://github.com/FenTechSolutions/CausalDiscoveryToolbox>.  
 1012  
 1013

### 1014 B.5 SORTNREGRESS

1015 SortnRegress algorithm includes Var-SortnRegress (Reisach et al., 2021) and  $R^2$ -  
 1016 SortnRegress (Reisach et al., 2023).  
 1017

1018 Reisach et al. (2021) emphasized that in the bivariate linear case, causal direction inferred by  
 1019 minimizing mean squared error (MSE) loss is from the variable with smaller variance to the variable  
 1020 with larger variance. They further hypothesized that, in the multivariate case, there is a consistency  
 1021 between the underlying causal direction of the data and the increasing order of the marginal variances  
 1022 of the variables. They provided a general definition of sortability:

$$1023 \mathbf{v}_\tau(X, \mathcal{G}) = \frac{\sum_{i=1}^d \sum_{(s \rightarrow t) \in A(\mathcal{G})^i} \text{incr}(\tau(X, s), \tau(X, t))}{\sum_{i=1}^d \sum_{(s \rightarrow t) \in A(\mathcal{G})^i} 1} \text{ where } \text{incr}(a, b) = \begin{cases} 1 & a < b \\ 1/2 & a = b \\ 0 & a > b \end{cases}, \quad (14)$$

$\tau$  represents a function with  $\tau(X) \in [0, 1]$ ,  $A(\mathcal{G})$  is the adjacency matrix of  $\mathcal{G}$ ,  $A(\mathcal{G})^i$  is the  $i$ -th matrix power,  $(s \rightarrow t) \in A(\mathcal{G})^i$  if and only if at least one directed path of length  $i$  from  $X_s$  to  $X_t$ .

In Var-Sortability,  $\tau(X, s)$  denotes the variance of  $X_s$ . Var-Sortability measures the consistency between the causal structure order and the increasing order of marginal variances of nodes. Intuitively, the greater Var-Sortability of the data, the better performance of methods based on MSE loss. Subsequently, they proposed the Var-SortnRegress algorithm to discover causality using only variance. In the first step, the nodes are ranked according to the increasing order of marginal variances. In the second step, linear and Lasso regression are used for estimation. We use the implementation of Var-SortnRegress algorithm in CausalDisco python package provided by the authors, available at <https://github.com/CausalDisco/CausalDisco>.

Reisach et al. (2021) demonstrated that the performance of linear differentiable causal discovery algorithms is greatly affected by Var-Sortability. Building on previous work, Reisach et al. (2023) pointed out that the coefficient of determination  $R^2$  remains unchanged after scaling the data, and proposed the  $R^2$ -SortnRegress algorithm, which achieves better performance on scale-variant data. The definition of  $R^2$ :

$$R^2 = 1 - \frac{\text{Var}(X_t - \mathbb{E}[X_t | X_{\{1, \dots, d\} \setminus \{t\}}])}{\text{Var}(X_t)}. \quad (15)$$

In  $R^2$ -Sortability,  $\tau(X, s)$  denotes the coefficient of determination  $R^2$  of  $X_s$ .  $R^2$ -Sortability measures the consistency between the causal structure order and the increasing order of  $R^2$ . If  $\mathbf{v}_{R^2}(X, \mathcal{G}) = 1$ , the causal order can be fully identified by the increasing order of  $R^2$ . If  $\mathbf{v}_{R^2}(X, \mathcal{G}) = 0$ , the causal order can be fully identified by the decreasing order of  $R^2$ . The only difference between  $R^2$ -SortnRegress and Var-SortnRegress lies in the definition of  $\tau$ . We use the implementation of  $R^2$ -SortnRegress algorithm in CausalDisco python package provided by the authors, available at <https://github.com/CausalDisco/CausalDisco>.

## B.6 NOTEARS

Based on (3), the NOTEARS score function is:

$$\min_{\mathcal{G}} F(\mathcal{G}; \mathbf{X}) = \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W(\mathcal{G})\|_F^2 + \lambda \|W(\mathcal{G})\|_1 \quad \text{s.t.} \quad h(W(\mathcal{G})) = 0, \quad (16)$$

where  $\|\cdot\|_F$  is the Frobenius norm,  $\|\cdot\|_1$  is the sum of absolute values of all elements in the matrix.

The unconstrained objective function obtained through the ALM is:

$$\min_{\mathcal{G}} L_{\mu}(W(\mathcal{G}), \theta, \alpha) = \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W(\mathcal{G})\|_F^2 + \lambda \|W(\mathcal{G})\|_1 + \alpha_t h(W(\mathcal{G})) + \frac{\mu_t}{2} |h(W(\mathcal{G}))|^2. \quad (17)$$

The update rule for the parameters is:

$$\begin{aligned} W(\mathcal{G})_k, \theta_k &= \arg \min_{W(\mathcal{G}), \theta} L_{\mu}(W(\mathcal{G}), \theta, \alpha) \\ \alpha_{k+1} &= \alpha_k + \mu_k h(W(\mathcal{G})_k) \\ \mu_{k+1} &= \begin{cases} \eta \mu_k, & \text{if } |h(W(\mathcal{G})_k)| > \gamma |h(W(\mathcal{G})_{k-1})| \\ \mu_k, & \text{otherwise} \end{cases}, \end{aligned} \quad (18)$$

where  $\theta$  represents the parameters of a neural network used to fit a nonlinear function, and  $\theta$  can be ignored for a linear model. The hyperparameters are usually set as  $\eta = 10$  and  $\gamma = \frac{1}{4}$ .

In practice, the optimization stopping criterion is  $h(W(\mathcal{G})_k) < \epsilon \in \{1e^{-6}, 1e^{-8}, 1e^{-10}\}$ , which does not guarantee the output to be a DAG. Finally, for values in  $W(\mathcal{G})$  with absolute values smaller than a threshold  $\tau$ , we set them to 0 in order to obtain a DAG as closely as possible. We use the implementation of NOTEARS algorithm provided by the authors, available at <https://github.com/xunzheng/notears>.

## B.7 GOLEM

GOLEM (Ng et al., 2020) proposed an improved loss function to address the numerical and ill-conditioned issues that often arise during the multiple iterations of optimization in NOTEARS. The

unconstrained optimization problem formulated by GOLEM is:

$$\min_{W(\mathcal{G}) \in \mathbb{R}^{d \times d}} \mathcal{S}_i(W(\mathcal{G}); \mathbf{X}) = \mathcal{L}_i(W(\mathcal{G}); \mathbf{X}) + \lambda_1 \|W(\mathcal{G})\|_1 + \lambda_2 h(W(\mathcal{G})), \quad (19)$$

where  $\lambda_1$  and  $\lambda_2$  are hyperparameters,  $i \in \{1, 2\}$ .

When assuming linear Gaussian with non-equal variances, that is, GOLEM-NV:

$$\mathcal{L}_1(W(\mathcal{G}); \mathbf{X}) = \frac{1}{2} \sum_{i=1}^d \log \left( \sum_{k=1}^n \left( X_i^{(k)} - W(\mathcal{G})_i^\top X^{(k)} \right)^2 \right) - \log |\det(I - W(\mathcal{G}))|, \quad (20)$$

where  $X_i^{(k)}$  denotes  $k$ -th data point of  $X_i$ .

When assuming linear Gaussian with equal variances, that is, GOLEM-EV:

$$\mathcal{L}_2(W(\mathcal{G}); \mathbf{X}) = \frac{d}{2} \log \left( \sum_{i=1}^d \sum_{k=1}^n \left( X_i^{(k)} - W(\mathcal{G})_i^\top X^{(k)} \right)^2 \right) - \log |\det(I - W(\mathcal{G}))|. \quad (21)$$

The authors proved that in the case of linear Gaussian with equal variances, when the hard DAG constraint is not satisfied, the least-squares optimal solution of NOTEARS returns a cyclic graph, whereas the optimal solution of GOLEM-EV corresponds to the ground-truth. GOLEM combines the maximum likelihood objective function with a soft DAG constraint, replacing the least-squares objective function and hard DAG constraint, making the optimization easier to solve and the results better. We use the implementation of GOLEM-EV algorithm provided by the authors, available at <https://github.com/ignavierng/golem>.

## B.8 NOTEARS-MLP

NOTEARS-MLP (Zheng et al., 2020) extends the differentiable causal discovery framework to the nonlinear case. Each variable  $X_j$  is defined as follows:

$$X_j = f_j(X_{pa(X_j)}, U_j), \forall j = 1, \dots, d. \quad (22)$$

The authors proved that  $f_j$  is independent of  $X_k$  if and only if  $\|\partial_k f_j\|_{L^2} = 0$ , where  $\|\cdot\|_{L^2}$  is the  $L^2$ -norm. Next, they define nonlinear causal effects through partial derivatives:

$$[W(f)]_{kj} = \|\partial_k f_j\|_{L^2}, \quad (23)$$

where  $[W(f)]_{kj}$  is the causal effects from  $X_k$  to  $X_j$ , and  $W(f) = W(f_1, \dots, f_d) \in \mathbb{R}^{d \times d}$ .

In practice, neural networks are used to fitting nonlinear functional relationships  $f_j$ :

$$\text{MLP}_j(\mathbf{X}; A^{(1)}, \dots, A^{(h)}) = \sigma \left( A^{(h)} \sigma \left( \dots A^{(2)} \sigma \left( A^{(1)} \mathbf{X} \right) \right) \right), \quad (24)$$

where  $A^{(\ell)} \in \mathbb{R}^{m_\ell \times m_{\ell-1}}$ ,  $\sigma$  is the activation function.

For the convenience of derivative computation, the authors proved that  $\text{MLP}_j$  are independent of  $X_k$  if and only if the  $k$ -th column of the first-layer weight matrix  $A^{(1)}$  is entirely zero. The parameters of  $\text{MLP}_j$  are  $\theta_j = (A_j^{(1)}, \dots, A_j^{(h)})$ . The authors ultimately obtain a weighted adjacency matrix representation that is independent of the depth of the neural network:

$$[W(\theta)]_{kj} = \left\| \text{th} - \text{column} \left( A_j^{(1)} \right) \right\|_2. \quad (25)$$

The objective function of NOTEARS-MLP is:

$$\min_{\theta} \frac{1}{n} \sum_{j=1}^d L(X_j, \text{MLP}_j(\mathbf{X}; \theta_j)) + \lambda \left\| A_j^{(1)} \right\|_1 \quad \text{s.t.} \quad h(W(\theta)) = 0. \quad (26)$$

NOTEARS-MLP trains  $d$  neural networks and represents the acyclicity constraint only through the first-layer parameters of neural networks. We use the implementation of NOTEARS-MLP algorithm provided by the authors, available at <https://github.com/xunzheng/notears>.

## 1134 B.9 GRAN-DAG

1135  
1136 GraN-DAG also extends NOTEARS to the nonlinear case and considers the ANM data generation  
1137 mechanism:

$$1138 X_j = f_j(X_{pa(X_j)}) + U_j, \forall j = 1, \dots, d. \quad (27)$$

1139 The authors state that the parameters of  $j$ -th neural network (NN) are:

$$1141 \phi_{(j)} = \left\{ W_{(j)}^{(1)}, \dots, W_{(j)}^{(L+1)} \right\}, \quad (28)$$

1142 where  $W_{(j)}^{(\ell)}$  is the  $\ell$ -th weight matrix of the  $j$ -th NN.

1143 They define the  $j$ -th connection matrix:

$$1144 C_{(j)} = \left| W_{(j)}^{(L+1)} \right| \dots \left| W_{(j)}^{(2)} \right| \left| W_{(j)}^{(1)} \right|. \quad (29)$$

1145 Based on the above definition, the authors construct a weighted adjacency matrix related to the depth  
1146 of the NN:

$$1147 (W_\phi)_{ij} = \begin{cases} \sum_{k=1}^m (C_{(j)})_{ki}, & \text{if } j \neq i \\ 0, & \text{otherwise} \end{cases}, \quad (30)$$

1148 where  $m$  is the output dimension of the NN.

1149 The objective function of GraN-DAG is:

$$1150 \max_{\phi} \mathbb{E}_{X \sim P(X)} \sum_{j=1}^d \log p_j(X_j | X_{pa(X_j)}; \phi_{(j)}) \quad \text{s.t. } h(\phi) = \text{Tr } e^{W_\phi} - d = 0. \quad (31)$$

1151 The unconstrained objective function of GraN-DAG is:

$$1152 \max_{\phi} \mathcal{L}(\phi, \alpha_t, \mu_t) = \mathbb{E}_{X \sim P(X)} \sum_{j=1}^d \log p_j(X_j | X_{pa(X_j)}; \phi_{(j)}) - \alpha_t h(\phi) - \frac{\mu_t}{2} h(\phi)^2. \quad (32)$$

1153 When the data generation mechanism follows the nonlinear Gaussian additive noise model, it can  
1154 be proven that the optimal solution of GraN-DAG corresponds to the ground-truth. We use the  
1155 implementation of GraN-DAG algorithm provided by the authors, available at <https://github.com/kurowasan/GraN-DAG>.

## 1156 B.10 NOCURL

1157 Since a DAG is related to curl-free functions on its edge set, NoCurl proposed a new representation  
1158 of a DAG:

$$1159 A = \gamma(W, p), \quad (33)$$

1160 where  $W$  is a skew-symmetric matrix with  $W = -W^T$ ,  $p \in \mathbb{R}^d$  is the potential function on the  
1161 vertices of the graph.

1162 The authors further proved that:

$$1163 \gamma(W, p) = W \circ \text{ReLU}(\text{grad}(p)), \quad (34)$$

1164 where  $\text{grad}$  is the gradient operator.

1165 The optimization problem established by NoCurl is:

$$1166 (W^*, p^*) = \underset{W, p}{\text{argmin}} F(\gamma(W, p), \mathbf{X}), \quad (35)$$

1167 with the optimal DAG  $A^* = W^* \circ \text{ReLU}(\text{grad}(p^*))$ . NoCurl implicitly ensures the acyclicity  
1168 constraint, overcoming the shortcomings of the ALM, avoiding multiple iterations, and improving  
1169 computational efficiency. We use the implementation of NoCurl algorithm provided by the authors,  
1170 available at <https://github.com/fishmoon1234/DAG-NoCurl>.

## B.11 DAGMA

DAGMA (Bello et al., 2022) proposed a log-determinant form of acyclicity representation  $h_{\text{ldet}}^s(W) = -\log \det(sI - W \circ W) + d \log s$  ( $s > 0$ ), which has three advantages compared to the exponential acyclicity constraints  $h_{\text{expm}}(W) = \text{Tr}(e^{W \circ W}) - d$  (Zheng et al., 2018) and polynomial acyclicity constraints  $h_{\text{poly}}(W) = \text{Tr}[(I + \alpha W \circ W)^d] - d$  ( $\alpha > 0$ ) (Yu et al., 2019). The authors proved that:

$$\begin{aligned} h_{\text{expm}}(W) &= \sum_{k=0}^{\infty} \frac{1}{k!} \text{Tr}((W \circ W)^k) - d, \\ h_{\text{poly}}(W) &= \sum_{k=0}^d \frac{\binom{d}{k}}{d^k} \text{Tr}((W \circ W)^k) - d, \end{aligned} \quad (36)$$

where  $\text{Tr}((W \circ W)^k)$  represents the information of cycles of length  $k$ . The information of cycles of length  $k$  in  $h_{\text{expm}}(W)$  and  $h_{\text{poly}}(W)$  is weakened by  $\frac{1}{k!}$  and  $\frac{\binom{d}{k}}{d^k}$ , respectively. It can be theoretically proven that  $h_{\text{ldet}}^s(W)$  is an upper bound of  $h_{\text{expm}}(W)$  and  $h_{\text{poly}}(W)$ , retaining more information about the cycles.

The authors also proved that:

$$\begin{aligned} \nabla h_{\text{expm}}(W) &= 2(e^{W \circ W})^\top \circ W \\ \nabla h_{\text{poly}}(W) &= 2 \left( \left( I + \frac{1}{d} W \circ W \right)^{d-1} \right)^\top \circ W. \\ \nabla h_{\text{ldet}}^s(W) &= 2((sI - W \circ W)^{-1})^\top \circ W \end{aligned} \quad (37)$$

$\nabla h_{\text{expm}}(W)$  and  $\nabla h_{\text{poly}}(W)$  are prone to the vanishing gradient problem. It can be theoretically proven that  $\nabla h_{\text{ldet}}^s(W)$  is an upper bound of  $\nabla h_{\text{expm}}(W)$  and  $\nabla h_{\text{poly}}(W)$ , retaining more information about the cycles.

The third advantage is that, in practice,  $h_{\text{ldet}}^s(W)$  and  $\nabla h_{\text{ldet}}^s(W)$  are faster to compute. Because the computation of  $h_{\text{ldet}}^s(W)$  and  $\nabla h_{\text{ldet}}^s(W)$  involves matrix log-determinant and matrix inverse, both of which have been extensively studied and solved. In contrast, other acyclicity constraints and their partial derivatives involve multiple matrix-to-matrix multiplications, which are slower. We use the implementation of DAGMA algorithm provided by the authors, available at <https://github.com/kevinsbello/dagma>.

## C RELATED WORK

**Differentiable causal discovery methods.** Building on traditional score-based causal discovery algorithms, NOTEARS (Zheng et al., 2018) transformed discrete constrained optimization into smooth equality-constrained optimization. This formulation has been extended to various settings, including more efficient linear models (GOLEM (Ng et al., 2020), NoCurl (Yu et al., 2021), NOFEARS (Wei et al., 2020), LEAST (Zhu et al., 2021)), neural networks (NOTEARS-MLP (Zheng et al., 2020), GraN-DAG (Lachapelle et al., 2019), DAGMA (Bello et al., 2022), DARING (He et al., 2021), CASTLE (Kyono et al., 2020)), generative adversarial networks (SAM (Kalinathan et al., 2022)), variational autoencoders (D-VAE (Zhang et al., 2019)), graph neural network (GAE (Ng et al., 2019), DAG-GNN (Yu et al., 2019)), federated learning (FedDAG (Gao et al., 2021)), reinforcement learning (RL-BIC (Zhu et al., 2019)), interventional data (DCDI (Brouillard et al., 2020)), time series data (DYNOTEARS (Pamfil et al., 2020)), multi-domain data (DICD (Wang et al., 2022), ReScore (Zhang et al., 2023), CASPER (Liu et al., 2023)), and domain adaptation (CAE (Yang et al., 2021)). Although differentiable causal discovery has made significant progress, it is also affected by Var-Sortability (Reisach et al., 2021; 2023) and highly non-convex optimization problems (Ng et al., 2024). [Recent research by Deng et al. \(2024\) shows that differentiable causal discovery methods can achieve scale invariance and global optimization when the correct loss function is used.](#)

## D TABLE RESULTS FOR RUNTIME OF THE BENCHMARK METHODS

The results in Table 8, 9, 10 and 11 show that differentiable causal discovery, exemplified by DAGMA, NOTEARS-MLP, and NoCurl, achieve superior performance with almost negligible runtime cost.

Table 8: Results for runtime (in seconds) on degree  $k = 2$  graphs of 10 and 20 nodes. The reported results are the mean and standard deviation of the runtime over 10 repetitions across different graph types, vanilla and model assumption violation scenarios.

Method	$d$	Runtime (seconds)
PC	10	1.29±0.24
	20	1.91±0.35
GES	10	1.64±0.53
	20	7.82±2.03
DirectLiNGAM	10	1.25±0.28
	20	2.06±0.43
Var-SortnRegress	10	1.11±0.25
	20	1.26±0.29
$R^2$ -SortnRegress	10	1.13±0.16
	20	1.24±0.35
NOTEARS	10	9.35±2.63
	20	35.57±4.71
GOLEM	10	130.23±2.54
	20	178.52±3.41
NoCurl	10	5.68±0.49
	20	10.29±1.84
CAM	10	45.73±2.67
	20	113.37±3.85
NOTEARS-MLP	10	4.70±0.82
	20	5.84±0.97
GraN-DAG	10	119.28±2.15
	20	211.59±4.35
DAGMA	10	<b>2.41±0.32</b>
	20	<b>3.19±0.48</b>



## E TABLE RESULTS ACROSS NODES, GRAPH TYPES, AND GRAPH DENSITIES

More experimental results for linear, nonlinear and MLP settings are reported in the Appendix as Table 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21 shows.

Table 9: Linear Setting, for ER-2 graphs of 10, 20, 50 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>
Random	25.6±3.1	57.9±8.5	27.9±2.3	67.8±7.8	25.9±3.5	60.4±11.3	27.9±3.2	62.0±8.1	26.3±3.5	57.7±7.6	26.1±3.7	60.7±12.9	26.7±3.0	64.0±8.7	27.5±3.4	63.1±6.0	28.3±1.0	63.6±7.8
PC	12.4±3.1	40.9±11.4	18.1±4.7	58.1±5.6	19.4±4.1	48.0±13.1	14.5±2.0	44.8±9.5	13.8±2.6	47.5±10.3	13.9±3.2	40.6±10.4	12.4±3.1	40.9±11.4	13.0±4.7	44.8±10.6	17.1±2.5	64.9±10.0
GES	13.8±7.8	32.0±15.6	25.9±7.7	42.6±16.0	20.2±4.8	46.2±16.7	20.8±5.5	49.7±11.5	15.5±6.1	35.2±12.2	17.6±6.4	39.0±15.7	13.8±7.8	32.0±15.6	10.1±5.2	25.6±12.6	16.2±2.2	57.8±10.7
DirectLiNGAM	19.6±3.3	46.1±10.6	20.4±5.0	42.0±16.0	17.6±2.4	48.8±12.4	19.7±3.2	50.4±8.4	16.3±3.9	34.7±9.9	19.7±4.3	44.7±11.9	21.8±4.3	62.6±9.3	20.1±4.3	49.8±11.1	20.0±0.0	63.5±7.7
Var-SomReg	11.2±3.5	8.4±8.5	17.6±5.8	12.6±9.9	19.6±2.8	11.4±8.7	18.8±2.4	16.5±10.6	17.9±3.3	8.6±9.3	12.4±3.1	13.5±8.0	30.7±5.1	34.5±10.3	7.3±3.5	9.3±8.4	15.6±3.3	39.0±6.7
I <sup>2</sup> -SomReg	20.2±4.8	32.4±14.0	25.7±4.4	37.6±15.0	25.6±2.0	39.2±16.0	25.6±4.9	38.8±19.6	26.0±5.2	37.0±14.4	28.6±4.8	51.0±11.3	20.2±4.8	32.4±14.0	20.5±6.7	32.8±8.3	20.3±3.7	66.1±9.7
NOTEARS	1.5±1.6	1.8±4.2	8.5±3.9	9.5±8.1	12.5±2.0	19.6±8.6	12.3±3.6	27.5±14.2	5.5±2.7	5.4±8.1	2.7±3.1	3.1±5.1	18.0±1.2	60.5±7.3	2.3±1.7	6.4±8.6	19.0±0.9	58.3±8.0
GoLEM	1.4±1.4	0.4±1.2	4.2±2.8	14.2±9.8	17.8±2.5	43.3±13.3	16.6±4.0	34.9±16.9	6.5±4.5	9.8±8.1	2.1±2.2	1.6±1.8	17.5±1.2	64.6±5.8	17.4±1.7	6.2±10.4	18.6±1.6	52.7±4.3
NoCut	2.0±1.8	5.1±5.8	9.1±4.2	5.4±3.9	11.8±1.8	17.9±8.4	14.6±2.5	17.5±10.8	6.6±2.9	5.5±5.7	2.2±2.3	2.0±4.4	27.2±5.1	69.9±7.9	3.1±3.2	4.7±5.8	19.1±1.0	58.9±9.5
DAGMA	1.2±1.2	3.3±5.3	8.8±7.7	12.6±2.5	18.5±8.6	12.3±3.6	28.4±15.3	5.5±2.3	12.0±8.2	2.1±2.2	0.6±1.8	17.9±1.4	58.7±6.8	15.1±1.4	4.5±7.4	19.0±0.9	58.3±8.0	

Table 10: Nonlinear Setting, for ER-2 graphs of 10, 20, 50 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>
Random	637.2±24.3	1393.2±62.2	652.6±21.3	1416.5±199.4	642.9±11.6	1456.3±96.8	635.5±15.6	1431.1±147.4	641.7±17.0	1387.0±118.3	640.7±18.0	1406.1±147.9	634.8±10.4	1456.3±112.2	632.7±16.3	1372.7±178.7	635.5±20.2	1413.6±154.1
PC	146.4±27.3	781.7±103.0	300.0±18.0	884.4±175.7	104.4±11.1	777.6±122.6	178.6±31.9	683.8±126.1	130.2±11.8	702.2±144.2	143.7±21.8	739.4±113.9	138.5±11.4	1248.0±112.0	158.3±19.3	810.2±120.6	99.8±6.7	1051.4±84.3
GES	106.0±16.3	113.1±16.8	93.0±16.1	145.7±19.3	281.6±27.1	208.3±51.2	427.3±48.8	207.3±88.8	309.5±48.3	121.8±43.7	138.5±27.7	127.6±54.4	66.5±41.0	115.4±154.3	112.6±31.3	86.1±62.3	86.4±6.2	85.4±7.0
DirectLiNGAM	388.2±7.2	555.8±142.2	406.3±104.6	401.3±104.6	346.0±138.3	903.4±138.6	552.5±56.8	853.3±129.0	549.4±68.9	573.7±133.0	442.7±109.4	642.1±176.9	388.3±72.4	555.8±142.2	420.2±68.4	547.0±78.6	108.1±7.9	133.4±89.6
Var-SomReg	15.2±6.2	71.4±5.5	46.5±15.4	266.5±9.3	91.0±18.2	610.0±88.3	122.4±37.7	94.2±134.6	23.5±4.4	181.3±15.5	16.7±4.5	66.6±84.9	90.8±3.8	972.0±8.0	91.2±9.6	92.4±1.9	95.7±0.7	
I <sup>2</sup> -SomReg	11.8±4.7	58.9±8.1	52.7±8.6	290.4±9.3	97.2±1.0	102.8±14.7	99.2±5.0	102.8±19.3	25.1±7.0	142.7±60.3	13.2±4.1	64.3±46.7	82.4±7.7	952.1±84.5	12.2±2.8	79.8±5.5	94.9±2.7	99.7±0.9
NOTEARS	22.8±12.1	33.2±10.6	15.6±2.1	42.4±8.8	20.0±1.1	63.8±11.3	16.2±2.3	44.2±10.0	14.8±2.1	40.5±10.8	11.2±3.1	37.7±10.0	10.0±3.7	26.1±12.0	10.4±3.4	25.7±8.1	16.3±2.1	54.7±5.1
GoLEM	12.4±4.3	47.5±39.7	42.7±10.1	284.2±107.8	81.2±9.4	635.0±21.0	103.0±12.3	93.2±119.9	24.4±7.6	170.7±62.0	14.6±4.4	55.0±33.3	90.6±4.0	969.3±88.7	13.3±4.3	92.4±1.9	99.0±0.7	
NoCut	13.5±2.0	40.7±8.1	18.6±2.0	62.0±13.3	17.3±1.3	54.9±8.3	19.0±2.0	56.6±10.5	16.4±4.3	36.6±15.7	8.6±3.1	29.6±14.8	15.7±2.7	53.3±12.9	13.7±2.7	41.4±7.8	3.3±3.1	9.5±11.6
DAGMA	13.5±2.0	40.7±8.1	18.6±2.0	62.0±13.3	17.3±1.3	54.9±8.3	19.0±2.0	56.6±10.5	16.4±4.3	36.6±15.7	8.6±3.1	29.6±14.8	15.7±2.7	53.3±12.9	13.7±2.7	41.4±7.8	3.3±3.1	9.5±11.6

Table 11: MLP Setting, for ER-2 graphs of 10, 20, 50 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>	SHD <sub>i</sub>	SHD <sub>j</sub>
Random	27.7±3.2	63.6±11.2	27.4±2.5	59.3±9.4	27.4±3.6	63.2±7.7	28.5±2.3	65.8±8.6	24.9±3.2	62.1±10.1	28.7±2.1	69.8±7.9	27.5±2.2	69.9±4.9	27.0±4.2	68.1±8.0	28.3±3.4	64.5±10.7
PC	17.1±2.5	64.9±10.0	18.5±1.9	62.8±9.1	18.2±1.1	61.6±12.2	18.4±1.3	61.2±9.7	21.3±3.2	56.5±10.4	15.4±1.4	57.3±9.5	17.1±2.5	64.9±10.0	17.6±3.1	63.5±10.5	12.4±3.1	40.9±13.4
GES	16.2±2.2	57.8±10.7	17.1±2.1	55.2±15.5	17.0±1.0	52.0±18.7	17.2±1.0	54.0±7.8	19.8±2.5	55.7±18.4	15.0±4.5	50.2±13.1	16.2±2.2	57.8±10.7	15.6±2.6	56.2±9.9	13.8±7.8	32.0±13.6
DirectLiNGAM	4.7±1.9	16.3±5.5	10.4±2.8	38.2±5.3	13.9±1.9	48.2±7.4	12.5±1.0	33.9±16.3	13.8±2.9	21.7±12.8	7.1±3.2	24.7±13.0	17.1±4.9	16.3±5.5	6.5±2.1	15.4±7.3	4.7±1.9	45.9±17.0
Var-SomReg	12.4±2.2	36.3±7.1	17.0±1.7	49.2±8.6	16.5±1.8	48.9±4.8	17.0±1.7	47.7±11.0	16.4±3.7	42.6±10.6	11.4±2.1	43.8±9.2	16.1±2.5	48.3±9.7	12.3±2.1	33.6±7.0	5.9±2.5	19.7±8.7
I <sup>2</sup> -SomReg	12.7±2.4	33.2±10.6	15.6±2.1	42.4±8.8	20.0±1.1	63.8±11.3	16.2±2.3	44.2±10.0	14.8±2.1	40.5±10.8	11.2±3.1	37.7±10.0	10.0±3.7	26.1±12.0	10.4±3.4	25.7±8.1	16.3±2.1	54.7±5.1
NOTEARS	13.5±2.0	40.7±8.1	18.6±2.0	62.0±13.3	17.3±1.3	54.9±8.3	19.0±2.0	56.6±10.5	16.4±4.3	36.6±15.7	8.6±3.1	29.6±14.8	15.7±2.7	53.3±12.9	13.7±2.7	41.4±7.8	3.3±3.1	9.5±11.6
GoLEM	10.7±7.8	251.0±36.8	102.5±8.9	237.5±33.7	103.5±3.2	230.9±25.5	105.6±2.7	249.0±31.2	101.5±7.1	256.4±27.0	103.5±6.1	232.6±21.6	106.7±8.8	247.7±36.9	102.7±7.8	233.2±29.6	101.1±5.7	234.6±42.6
NoCut	34.5±4.4	226.3±95.7	37.8±1.8	202.0±22.8	36.2±2.0	213.3±13.9	39.6±6.3	228.8±35.4	49.9±6.3	213.9±24.8	32.7±3.1	221.6±28.5	34.5±4.4	226.3±95.7	34.2±2.6	217.0±31.1	34.2±2.6	217.0±31.1
DAGMA	32.6±3.3	190.1±22.6	37.7±2.5	199.2±20.2	34.9±2.4	191.0±41.3	36.9±2.3	190.0±28.1	51.0±3.4	211.7±23.6	32.7±3.0	187.9±35.8	32.6±3.3	190.1±22.6	31.6±3.9	189.9±31.4	34.3±2.6	108.5±11.1
Var-SomReg	15.8±4.4	70.8±17.7	34.5±2.8	166.4±24.5	31.0±2.1	178.9±33.0	27.9±6.6	125.4±28.8	41.3±5.2	90.3±25.8	16.5±4.2	119.4±40.6	15.8±4.4	70.8±17.7	15.1±2.6	85.1±21.8	41.2±8.1	213.4±49.0
I <sup>2</sup> -SomReg	27.8±2.4	133.3±24.4	37.8±1.5	204.4±24.8	33.0±1.2	174.1±33.6	39.9±0.3	215.2±26.4	40.0±0.0	216.0±24.9	24.8±2.1	163.0±26.5	32.0±2.4	181.0±36.1	27.0±3.4	130.8±25.1	19.2±4.9	87.3±17.7
NOTEARS	30.8±4.2	148.5±26.6	48.1±2.1	218.2±29.3	41.2±7.0	208.0±24.4	36.8±2.9	194.3±22.2	42.4±1.1	161.7±30.9	28.3±4.0	162.2±30.6	31.0±4.3	150.7±26.2	31.0±3.2	149.4±27.1	37.7±2.5	189.3±26.5
DAGMA	27.3±4.4	141.9±36.4	36.2±4.3	175.0±21.9	35.1±2.1	186.3±15.5	39.9±1.1	213.5±27.5	40.0±0.0	216.0±24.9	19.8±6.0	125.1±33.0	32.2±2.6	193.4±27.6	37.3±3.3	141.2±39.0	14.8±5.8	69.7±35.0

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

Table 12: Linear Setting, for ER-4 graphs of 10, 20, 50 nodes.

Table with 10 columns: 10 nodes, 20 nodes, 50 nodes. Each column contains sub-tables for Vanilla model, Latent confounders, Measurement error, Autoregressive, Heterogeneous, Unfaithful, Scale-variant, Missing, and Mechanism violation. Each sub-table lists model names (Random, DirectINGAM, etc.) and their corresponding SHD, SIDI, and SDD values.

Table 13: Nonlinear Setting, for ER-4 graphs of 10, 20, 50 nodes.

Table with 10 columns: 10 nodes, 20 nodes, 50 nodes. Each column contains sub-tables for Vanilla model, Latent confounders, Measurement error, Autoregressive, Heterogeneous, Unfaithful, Scale-variant, Missing, and Mechanism violation. Each sub-table lists model names (Random, DirectINGAM, etc.) and their corresponding SHD, SIDI, and SDD values.

Table 14: Linear Setting, for SF-2 graphs of 10, 20, 50 nodes.

Table with 10 columns: 10 nodes, 20 nodes, 50 nodes. Each column contains sub-tables for Vanilla model, Latent confounders, Measurement error, Autoregressive, Heterogeneous, Unfaithful, Scale-variant, Missing, and Mechanism violation. Each sub-table lists model names (Random, DirectINGAM, etc.) and their corresponding SHD, SIDI, and SDD values.

Table 15: Nonlinear Setting, for SF-2 graphs of 10, 20, 50 nodes.

Table with 10 columns: 10 nodes, 20 nodes, 50 nodes. Each column contains sub-tables for Vanilla model, Latent confounders, Measurement error, Autoregressive, Heterogeneous, Unfaithful, Scale-variant, Missing, and Mechanism violation. Each sub-table lists model names (Random, DirectINGAM, etc.) and their corresponding SHD, SIDI, and SDD values.

Table 16: Linear Setting, for SF-4 graphs of 10, 20, 50 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>
Random	28.72±2	77.2±4	30.6±3	71.4±8	30.6±3	71.4±8	28.3±1	72.2±1	28.9±2	76.4±6	31.1±3	76.6±3	30.8±0	73.1±1	30.0±4	75.2±5	27.0±8	70.6±3
PC	27.3±2	71.4±1	26.1±2	76.6±2	25.8±1	70.8±0	27.1±1	73.2±5	28.2±4	75.7±2	26.2±3	74.2±4	27.3±2	71.4±1	26.6±3	71.6±2	26.3±7	77.7±6
GES	23.9±5	62.7±9	27.2±6	64.3±1	24.5±7	60.0±14	24.5±7	65.5±10	27.4±5	64.7±13	24.6±5	67.5±6	23.0±1	62.7±9	26.1±6	67.0±7	26.1±2	82.6±0
DirectLiNGAM	31.6±5	63.8±16	24.9±2	65.1±5	26.1±2	67.8±3	21.4±1	61.6±7	23.0±2	68.6±9	19.8±7	67.4±8	24.6±5	73.1±4	23.3±9	65.0±5	30.0±0	81.9±7
Var-SortReg	10.7±2	22.5±9	13.0±7	21.4±0	17.8±1	27.6±0	15.7±6	18.6±13	15.5±3	20.1±7	13.2±2	24.3±0	26.2±4	61.8±6	9.9±0	17.2±1	20.6±5	57.9±7
R <sup>2</sup> -SortReg	25.1±5	60.2±8	24.5±6	51.1±1	27.9±3	67.0±3	23.5±5	69.2±15	25.8±4	57.4±16	26.5±9	57.3±4	25.3±0	60.2±8	21.5±4	47.7±1	29.9±1	83.9±6
NOTEARS	3.2±0	16.2±4	8.0±5	18.5±1	15.3±5	35.3±12	10.4±8	21.6±12	10.9±8	27.1±19	2.4±1	14.7±11	24.7±9	29.2±6	4.0±7	18.8±0	29.1±7	77.7±2
GOLEM	2.2±1	9.0±8	9.3±7	21.3±4	24.6±2	69.4±2	18.4±1	44.1±5	9.1±3	21.5±6	1.6±1	7.0±7	25.7±5	79.9±1	1.4±1	6.9±7	28.2±1	73.1±7
NCut	4.0±2	14.0±1	10.4±2	21.9±1	15.5±7	60.0±1	14.1±7	32.7±2	10.2±6	21.3±5	2.5±1	17.0±6	26.8±2	59.9±2	2.6±2	11.7±6	29.0±1	77.7±6
DAGMA	2.0±2	10.9±1	8.5±2	20.6±5	15.2±3	37.9±7	15.5±8	25.7±10	6.3±5	20.5±10	1.3±1	7.3±16	25.7±2	79.3±4	1.1±1	6.4±6	29.1±7	77.7±2
20 nodes	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>
Random	113.4±6	314.7±10	109.9±5	305.6±7	110.9±7	308.7±10	109.9±4	303.2±11	114.5±7	310.0±14	113.3±9	309.1±10	113.1±7	301.9±8	111.3±7	303.0±15	115.3±8	311.8±6
PC	77.5±2	345.1±6	78.9±4	349.5±7	92.6±1	329.8±14	82.8±5	341.9±7	78.9±4	347.3±6	76.2±7	347.6±9	77.5±2	345.1±6	78.1±4	350.5±8	67.6±7	371.1±9
GES	96.6±4	257.5±25	131.0±1	235.8±3	110.0±8	293.5±25	115.2±19	241.3±20	127.5±14	259.1±23	89.9±0	250.7±17	96.6±4	257.5±25	83.2±0	267.3±29	103.2±0	368.9±2
DirectLiNGAM	58.8±10	283.4±7	76.9±6	280.1±1	67.1±8	318.0±12	81.3±4	284.8±7	72.9±8	292.2±8	65.6±8	296.3±6	68.3±9	344.1±1	60.2±1	287.6±20	69.9±3	368.5±2
Var-SortReg	54.7±0	104.0±23	105.0±3	88.9±5	111.2±7	119.4±17	108.9±3	107.1±3	106.3±9	100.0±3	106.8±9	118.2±3	108.5±7	215.3±2	57.9±9	108.5±12	61.4±8	109.1±1
R <sup>2</sup> -SortReg	128.5±0	251.4±6	137.9±9	229.1±25	137.4±7	243.4±6	133.3±9	225.8±14	143.7±3	246.9±12	137.5±1	261.9±1	128.5±0	251.4±6	127.5±9	271.4±7	151.2±1	371.7±1
NOTEARS	12.5±8	76.2±6	41.1±9	109.1±24	57.5±8	207.7±3	84.5±6	185.3±0	30.6±7	160.3±4	10.3±6	68.6±2	68.2±5	35.1±5	10.5±8	87.8±1	68.8±6	349.0±2
GOLEM	8.5±4	83.9±0	52.7±4	136.2±7	69.8±4	360.6±1	76.5±8	365.9±4	40.1±8	190.1±1	10.1±5	82.3±9	61.6±2	35.8±0	7.9±0	72.7±8	67.9±9	346.3±8
NCut	10.7±6	64.9±7	55.8±5	125.3±3	55.6±9	236.5±8	89.3±2	160.3±8	43.3±8	120.6±7	14.5±1	76.5±0	73.0±7	336.5±2	9.7±1	63.1±9	68.8±7	350.8±6
DAGMA	83.6±2	58.1±8	41.6±5	111.2±3	55.6±7	270.6±14	77.3±19	207.9±3	35.9±7	137.0±1	7.2±3	38.9±2	65.6±4	345.2±0	7.5±3	52.5±6	68.8±6	349.0±2
50 nodes	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>
Random	659.9±24	1901.2±46	670.5±28	1899.4±20	658.5±18	1904.2±1	667.6±18	1908.4±7	668.3±19	1914.0±29	659.2±28	1918.7±29	664.4±18	1914.4±13	669.1±11	1911.2±24	648.4±24	1911.2±27
PC	2071.2±8	2126.2±10	2173.2±1	1979.3±0	2618.1±1	2323.8±1	317.5±17	1807.5±12	228.4±2	2095.8±2	2106.2±5	2186.19	226.4±8	2347.1±19	250.3±20	2100.7±9	1900.4±0	2405.9±1
GES	3471.1±9	622.1±0	6807.1±3	694.0±1	6954.9±1	1122.8±0	70.7±2	674.8±4	74.3±6	611.0±2	441.3±4	800.0±1	660.6±1	1515.1±8	352.9±7	617.0±8	1700.7±4	2284.0±9
DirectLiNGAM	777.4±5	163.7±6	968.6±17	1489.3±8	788.3±9	1721.3±5	914.1±7	1419.7±1	937.2±6	1530.2±4	902.2±9	1639.2±9	791.4±6	1637.5±6	791.4±6	1637.5±6	206.2±5	247.0±1
Var-SortReg	62.6±1	984.5±2	122.7±1	182.4±7	205.6±2	2143.9±1	324.0±7	221.9±5	153.1±3	158.1±2	60.2±8	952.6±2	181.3±7	235.6±7	61.6±6	111.8±6	188.7±9	237.0±1
R <sup>2</sup> -SortReg	65.1±9	1062.9±1	193.1±7	324.3±6	106.6±2	231.6±1	194.4±9	299.1±1	152.4±3	181.1±2	55.0±0	117.1±9	60.9±1	289.4±1	66.0±1	138.4±6	188.0±9	409.8±2
NOTEARS	24.6±0	254.9±2	34.8±3	89.1±0	213.8±1	2045.1±8	44.5±0	1289.7±4	175.9±2	872.5±2	36.1±8	374.9±2	241.1±2	271.2±5	31.4±4	291.5±1	188.4±9	238.0±2
GOLEM	21.9±4	48.9±7	113.3±4	91.7±1	197.2±4	2387.8±6	237.5±4	2326.0±7	172.0±2	1376.2±6	33.3±2	396.7±9	181.6±9	237.7±0	24.1±4	42.5±1	188.7±9	275.7±2
NCut	24.6±0	254.9±2	34.8±3	89.1±0	213.8±1	2045.1±8	44.5±0	1289.7±4	175.9±2	872.5±2	36.1±8	374.9±2	241.1±2	271.2±5	31.4±4	291.5±1	188.4±9	238.0±2
DAGMA	21.9±4	48.9±7	113.3±4	91.7±1	197.2±4	2387.8±6	237.5±4	2326.0±7	172.0±2	1376.2±6	33.3±2	396.7±9	181.6±9	237.7±0	24.1±4	42.5±1	188.7±9	275.7±2

Table 17: Nonlinear Setting, for SF-4 graphs of 10, 20, 50 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>
Random	29.5±2	77.8±1	29.9±0	75.6±1	28.2±7	73.9±5	31.6±2	78.0±1	32.1±0	76.4±7	29.2±7	73.9±4	30.3±4	72.7±5	28.5±2	71.9±5	30.0±1	74.9±5
PC	26.3±7	77.6±9	24.1±4	85.3±4	27.5±2	85.3±1	29.1±0	78.7±7	29.1±0	72.6±0	26.1±9	69.6±8	26.3±7	77.6±9	25.6±2	77.6±9	27.1±2	71.6±1
GES	26.1±2	82.6±9	26.1±5	78.4±6	28.3±1	82.9±5	28.2±2	71.1±8	26.8±2	71.1±8	21.3±1	60.5±1	26.1±2	82.6±9	25.0±2	79.8±5	30.6±0	62.1±8
CM	5.2±4	16.1±4	13.8±4	43.7±8	20.8±8	64.9±14	15.8±3	46.0±13	17.9±3	31.9±13	9.9±1	44.9±1	5.2±4	16.1±4	4.5±1	18.4±7	7.7±2	73.2±0
NOTEARS	20.1±4	99.0±2	27.0±4	75.7±7	31.7±4	121.1±1	26.9±8	74.1±4	21.2±7	88.3±7	15.2±7	88.3±7	21.2±7	91.0±0	12.1±9	91.0±0	12.1±9	46.1±7
NCut	10.7±4	12.4±7	20.3±9	51.7±2	27.3±1	75.9±1	23.8±4	54.8±13	19.4±3	40.9±13	15.9±3	48.9±19	11.5±6	18.6±1	12.2±6	40.1±7	17.9±2	65.1±2
DAGMA	22.9±2	62.7±9	23.9±2	73.2±5	28.0±7	73.2±5	29.7±5	10.1±5	22.5±5	56.7±10	11.5±2	40.1±7	20.2±7	81.6±2	7.2±4	44.3±8	17.4±7	34.9±9
20 nodes	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>	SHD <sub>i</sub>	SHD <sub>o</sub>
Random	111.2±5	307.7±1	114.4±8	302.9±0	111.6±7	303.8±7	116.0±9	311.1±13	114.4±7	306.5±10	110.6±5	303.1±9	114.2±8	307.5±6	116.0±7	305.3±8	109.1±6	305.7±9
PC	67.6±7	371.1±9	68.2±1	357.6±9	68.5±1	373.2±8	70.0±1	369.4±13	82.2±5	349.9±3	68.4±5	348.1±2	67.6±7	371.1±9	66.3±0	370.3±6	77.5±2	345.1±6
GES	62.2±0	369.9±2	67.1±2	358.0±7	65.9±1	362.5±3	68.1±2	362.9±5	81.4±5	357.4±10	70.9±5	294.3±8	62.3±0	369.9±2	60.0±3	369.9±2	79.3±9	268.2±2
DirectLiNGAM	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9	185.3±3	158.7±9
Var-SortReg	56.1±1	324.3±6	68.7±9	379.5±4	64.4±2	342.7±1	70.0±0	361.0±0	67.7±5	330.1±8	37.3±9	299.8±9	61.0±2	340.7±2	54.9±7	322.1±0	34.2±4	239.6±2
R <sup>2</sup> -SortReg	59.4±1	364.2±7	67.1±2	354.6±6	69.0±5	354.9±1	69.4±9	339.9±1	69.4±9	339.9±1	69.4±9	339.9±1	69.4±9	339.9±1	69.4±9	339.9±1	69.4±9	339.9±1
NOTEARS	58.3±3	349.3±1	65.3±1	358.3±9	69.9±5	360.1±0	67.6±7	324.7±2	35.8±3	278.4±2	35.8±3	278.4±2	60.0±2	372.4±1	57.4±1			

## F TABLE RESULTS FOR COMBINED MISSPECIFIED SCENARIOS

We consider two (confounding and heterogeneity), three (confounding, measurement error, and heterogeneity) and four (confounding, measurement error, heterogeneity, and autoregression) combined misspecified scenarios. The results in Table 22 and Table 23 show that under combined misspecified scenarios, the performance of various methods is worse compared to single misspecified scenario. However, differentiable causal discovery still achieves optimal or competitive performance.

Table 22: Linear Setting with two (confounding and heterogeneity), three (confounding, measurement error, and heterogeneity) and four (confounding, measurement error, heterogeneity, and autoregression) combined misspecified scenarios, for ER-2 graphs of 10 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Heterogeneous		Autoregressive		Two combined scenarios		Three combined scenarios		Four combined scenarios	
	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>
Random	29.1±1.1	68.4±2.7	25.1±4.7	65.2±5.7	24.3±3.9	65.6±6.6	26.2±4.3	70.4±2.3	27.5±3.6	65.8±4.2	27.1±1.3	63.4±4.5	25.7±4.2	61.2±7.2	28.1±2.3	69.2±3.5
PC	12.4±3.1	40.9±13.4	18.1±4.7	58.1±15.6	19.4±4.1	48.0±13.1	13.8±2.6	47.5±10.3	14.5±2.0	44.8±9.5	19.5±3.2	61.0±8.3	21.8±4.4	55.3±13.5	21.6±4.0	59.3±10.0
GES	13.8±7.8	32.0±15.6	25.9±7.7	42.6±14.0	20.2±4.8	46.2±16.7	15.5±6.1	35.2±12.2	20.8±5.5	49.7±11.5	28.1±3.4	41.2±12.8	27.4±1.7	62.0±8.9	24.7±5.0	50.2±11.1
DirectLINGAM	19.6±3.3	46.1±10.6	20.4±5.0	42.0±6.0	17.6±2.4	48.8±12.4	16.3±3.9	34.7±9.9	19.7±4.2	50.4±8.4	22.4±2.4	46.7±12.7	18.0±3.1	43.7±6.3	18.9±3.2	38.9±4.8
Var-SortnRegress	11.2±3.5	8.4±8.5	17.6±5.8	12.6±9.9	19.6±2.8	<b>11.4±8.7</b>	17.9±3.3	8.6±9.3	18.8±2.4	<b>16.5±10.6</b>	21.5±3.9	11.5±10.8	20.7±3.3	<b>13.0±9.9</b>	20.3±3.3	23.7±9.9
R <sup>2</sup> -SortnRegress	20.2±4.8	32.4±14.0	25.7±4.1	37.6±13.0	25.6±6.0	39.2±16.0	26.0±5.4	37.0±14.4	25.6±4.9	38.8±19.0	25.5±3.6	31.9±17.6	26.8±3.8	48.3±10.4	26.8±5.2	48.1±12.1
NOTEARS	1.5±1.6	1.8±4.2	8.5±3.9	9.5±8.1	12.5±2.0	19.6±8.6	<b>5.5±2.7</b>	<b>5.4±5.1</b>	<b>12.2±3.6</b>	27.5±14.2	12.3±3.4	33.1±5.7	<b>14.1±3.5</b>	25.2±6.8	18.6±3.9	31.8±10.1
GOLEM	1.4±1.4	<b>0.4±1.2</b>	<b>6.7±2.8</b>	14.2±9.8	17.8±2.5	43.1±13.3	6.5±4.5	9.8±8.1	16.6±4.0	34.9±16.9	<b>10.3±1.7</b>	24.3±5.0	18.6±3.7	46.2±9.1	19.7±0.5	72.7±6.3
NoCurl	2.0±1.8	5.1±5.8	9.1±4.2	<b>5.4±3.9</b>	<b>11.8±1.8</b>	<b>17.9±3.4</b>	6.6±2.9	5.5±5.7	14.8±2.5	<b>17.5±10.8</b>	13.4±3.4	<b>7.4±6.9</b>	15.4±3.3	<b>19.8±7.0</b>	20.2±3.5	<b>21.3±8.2</b>
DAGMA	<b>1.2±1.2</b>	3.3±3.5	8.4±3.9	8.8±7.7	12.6±2.5	18.5±8.6	<b>5.5±2.3</b>	12.0±8.2	<b>12.2±3.6</b>	28.4±15.3	12.0±3.8	12.3±8.9	14.9±4.4	20.4±6.6	<b>17.9±3.1</b>	50.2±14.2

Table 23: MLP Setting with two (confounding and heterogeneity), three (confounding, measurement error, and heterogeneity) and four (confounding, measurement error, heterogeneity, and autoregression) combined misspecified scenarios, for ER-2 graphs of 10 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Heterogeneous		Autoregressive		Two combined scenarios		Three combined scenarios		Four combined scenarios	
	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>
Random	27.5±2.6	62.6±7.3	29.4±1.1	62.8±9.4	29.3±1.2	69.4±3.1	28.9±1.8	70.4±2.8	27.9±2.3	72.2±5.5	28.4±2.0	67.2±4.5	24.2±4.1	66.2±6.3	29.6±2.3	68.5±5.6
PC	16.7±3.2	32.0±9.5	18.3±3.5	54.1±9.3	17.9±5.0	47.4±15.5	20.7±3.5	53.9±13.1	16.2±4.7	49.2±13.1	22.9±2.6	63.7±9.6	24.7±4.2	64.4±11.5	22.6±3.2	65.1±14.3
GES	21.9±4.4	46.9±8.6	28.2±7.5	53.1±18.6	20.2±5.4	54.7±14.3	27.2±4.9	48.6±10.9	18.9±5.2	43.2±15.5	28.6±2.9	53.0±8.4	28.5±3.5	62.9±6.3	24.6±2.3	61.6±6.7
CAM	12.4±3.6	39.3±16.5	16.6±4.2	42.4±17.7	19.7±4.7	59.6±12.0	19.9±3.6	49.1±7.2	16.0±3.4	46.0±16.6	23.4±3.1	60.9±6.1	25.4±3.3	53.8±13.4	24.4±1.9	63.4±8.8
NOTEARS-MLP	8.1±2.7	22.2±10.6	11.7±5.5	33.3±17.0	18.5±3.7	47.1±11.9	14.5±2.2	34.3±10.4	15.7±4.3	39.8±9.6	18.5±2.3	57.4±11.7	19.6±0.7	63.4±10.2	22.5±3.3	61.0±9.2
GrN-DAG	13.3±3.6	32.9±12.3	12.9±3.2	40.2±12.7	16.6±2.0	45.0±10.5	14.3±2.5	38.5±7.7	15.5±2.3	46.0±9.5	16.0±2.5	52.9±11.2	20.4±1.6	63.8±7.2	20.3±2.8	62.2±10.0
DAGMA	<b>6.2±1.7</b>	<b>18.2±8.7</b>	<b>9.3±4.3</b>	<b>27.8±10.8</b>	<b>14.1±2.6</b>	<b>39.2±8.7</b>	<b>12.7±2.9</b>	<b>31.6±7.2</b>	<b>13.6±2.1</b>	<b>41.2±3.6</b>	<b>15.3±2.7</b>	<b>48.2±10.4</b>	<b>18.4±0.9</b>	<b>51.6±10.4</b>	<b>19.6±2.5</b>	<b>60.8±7.3</b>

## G TABLE RESULTS FOR NON-GAUSSIAN NOISE

We consider the vanilla model with exponential noise. The results in Table 24 and Table 25 show that differentiable causal discovery still achieve optimal or competitive performance when model assumptions are violated.

Table 24: Linear Setting with exponential noise, for ER-2 graphs of 10 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>
Random	29.2±1.4	64.0±8.2	27.5±1.9	69.2±3.0	25.7±2.6	65.3±4.8	25.1±2.8	65.7±8.5	26.5±3.0	64.2±7.6	24.7±3.1	71.0±2.4	28.1±1.6	65.7±5.4	26.6±3.4	66.1±5.4	27.9±3.0	67.2±7.2
PC	11.7±2.3	38.7±7.3	20.1±2.6	58.6±8.7	20.5±3.3	46.8±15.0	17.9±2.2	50.2±8.8	20.7±4.3	58.4±14.5	14.0±3.1	44.6±13.5	11.7±2.3	38.7±7.3	14.4±4.8	42.2±15.3	17.9±3.6	54.4±12.2
GES	19.4±8.1	41.1±19.7	24.8±6.8	41.6±16.0	24.4±4.4	56.7±14.5	22.3±3.8	37.0±11.9	26.4±5.0	45.3±17.3	18.3±5.8	38.2±14.5	18.9±8.3	37.9±20.1	17.3±6.5	34.3±19.7	19.6±4.9	42.5±12.9
DirectLINGAM	<b>0.0±0.0</b>	<b>0.0±0.0</b>	13.6±5.9	20.1±15.4	15.0±3.7	27.1±9.5	12.0±3.4	18.1±12.2	9.1±4.3	10.3±9.2	1.4±2.2	3.9±4.9	<b>4.9±2.5</b>	<b>18.3±9.0</b>	<b>0.0±0.0</b>	<b>0.0±0.0</b>	16.5±3.5	45.9±10.1
Var-SortnRegress	6.8±4.5	9.7±7.8	17.0±5.2	10.0±6.6	18.4±4.0	17.2±9.1	16.2±2.7	10.2±7.9	22.1±2.1	10.9±6.9	8.7±6.1	15.0±12.0	26.9±5.5	58.2±8.3	5.7±3.8	10.6±11.3	10.4±5.0	28.8±7.9
R <sup>2</sup> -SortnRegress	15.7±6.7	30.4±12.2	27.9±4.1	35.9±16.1	26.4±5.9	44.5±14.9	23.1±5.0	37.3±13.9	26.6±3.4	32.6±14.3	20.9±6.3	34.3±15.0	15.7±6.7	30.4±12.2	15.8±7.1	27.0±10.0	29.0±4.2	55.8±11.0
NOTEARS	1.5±1.7	6.7±7.7	10.7±3.7	14.0±7.2	13.4±3.5	15.8±6.7	14.4±5.2	18.5±12.1	<b>6.6±2.7</b>	13.2±8.3	<b>0.1±0.3</b>	<b>0.7±2.1</b>	18.9±1.5	<b>65.3±8.7</b>	0.9±1.4	5.1±7.0	<b>14.1±4.6</b>	29.0±6.2
GOLEM	2.3±2.7	6.6±8.4	10.5±3.4	13.9±7.1	14.3±4.1	18.0±7.3	13.6±5.5	20.1±13.3	6.7±2.7	13.5±8.6	0.2±0.4	0.8±2.7	18.6±2.2	67.7±10.6	0.6±1.7	4.5±6.6	14.5±4.4	28.9±5.7
NoCurl	5.4±3.6	11.5±8.5	14.6±4.0	14.4±8.7	16.7±4.6	<b>14.3±8.6</b>	18.1±6.2	<b>16.8±13.7</b>	11.9±4.7	9.0±8.7	8.6±4.3	9.2±12.1	24.6±6.1	65.6±13.9	8.0±2.7	7.1±7.7	22.2±4.4	27.9±8.5
DAGMA	<b>0.0±0.0</b>	<b>0.0±0.0</b>	<b>10.1±2.9</b>	<b>13.8±6.9</b>	<b>13.2±2.4</b>	24.6±9.3	<b>12.9±4.7</b>	31.0±10.8	6.8±2.9	<b>4.7±5.3</b>	<b>0.1±0.3</b>	<b>0.7±2.1</b>	<b>18.2±4.3</b>	67.5±10.4	<b>0.3±0.9</b>	<b>0.6±1.8</b>	14.6±4.7	<b>28.8±7.3</b>

Table 25: MLP Setting with exponential noise, for ER-2 graphs of 10 nodes.

10 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>	SHD <sub>↓</sub>	SID <sub>↓</sub>
Random	27.9±2.2	62.6±6.6	26.3±2.5	62.6±5.9	28.0±1.3	66.9±6.9	25.0±3.2	68.1±3.2	29.8±1.1	67.6±4.3	26.1±4.0	60.3±5.7	30.0±1.0	69.9±2.6	25.4±5.0	63.2±8.4	26.7±4.8	64.7±9.2
PC	17.9±3.6	54.4±12.2	17.2±4.7	57.3±14.1	18.4±5.1	49.2±9.7	17.0±3.9	55.5±11.5	23.2±2.6	62.5±11.7	17.9±2.4	48.2±13.7	17.9±3.6	54.4±12.2	17.0±3.9	52.2±10.2	11.7±2.3	38.7±7.3
GES	19.6±4.9	42.5±12.9	22.8±7.0	42.8±18.0	19.6±4.7	46.6±12.4	17.1±4.8	47.0±16.2	30.0±2.3	58.5±8.4	23.6±6.0	44.1±12.3	19.6±4.9	42.2±12.5	20.7±4.3	44.2±14.5	19.4±8.1	41.1±9.7
CAM	11.2±2.6	40.3±17.3	14.1±4.9	35.5±14.1	21.2±3.8	61.6±11.2	13.7±3.5	42.3±16.4	23.3±3.5	63.5±12.6	12.5±5.2	35.7±21.4	<b>11.2±2.6</b>	<b>40.3±17.3</b>	12.3±2.1	41.6±20.6	21.2±7.6	62.8±13.1
NOTEARS-MLP	6.6±2.3	17.1±11.0	13.6±5.3	38.7±12.5	16.2±5.1	31.8±11.2	15.0±5.7	34.3±17.7	12.6±3.2	25.6±10.5	11.2±2.4	21.6±8.5	17.2±1.7	61.7±9.6	6.5±2.5	16.5±11.9	10.4±4.2	43.9±17.5
GrN-DAG	12.2±3.0	38.6±15.1	15.1±2.9	51.5±10.1	16.4±1.7	47.3±7.2	14.9±2.8	48.1±12.8	13.6±3.9	33.8±15.1	14.1±2.0	42.0±9.6	17.3±1.9	54.6±12.4	12.6±3.1	40.5±13.1	14.7±3.7	47.4±12.3
DAGMA	<b>4.5±2.5</b>	<b>13.5±9.3</b>	<b>12.1±5.2</b>	<b>32.0±9.7</b>	<b>14.3±5.0</b>	<b>28.2±9.3</b>	<b>12.2±5.6</b>	<b>28.5±12.5</b>	<b>10.8±3.0</b>	<b>21.8±8.6</b>	<b>9.2±1.7</b>	<b>19.5±7.6</b>	15.5±1.6	<b>52.0±9.7</b>	<b>5.3±2.4</b>	<b>15.2±11.8</b>	<b>9.9±3.1</b>	<b>38.3±14.6</b>

1512 H SUMMARY OF THE MOST COMPETITIVE METHODS

1513

1514

1515 Table 26: Summary of performances of the most competitive methods under linear setting. The  
 1516 reported results are the mean and standard deviation of the metrics over 10 repetitions across different  
 1517 graph types, vanilla and model assumption violation scenarios.

1518

1519

Method	$d$	SHD	SID
NOTEARS	10	8.51±5.92	23.88±23.32
	20	22.66±14.27	129.96±109.48
	50	61.07±39.27	945.01±774.12
GOLEM	10	9.02±6.90	30.67±27.39
	20	21.79±14.19	155.60±126.41
	50	55.57±35.00	1000.95±842.80
NoCurl	10	9.39±7.06	23.59±23.92
	20	32.22±23.51	99.43±99.13
	50	128.97±113.22	914.22±699.51
DAGMA	10	<b>8.17±6.12</b>	<b>22.42±22.72</b>
	20	<b>20.98±14.10</b>	<b>127.31±112.55</b>
	50	<b>55.19±37.87</b>	<b>882.76±853.40</b>

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533 Table 27: Summary of performances of the most competitive methods under nonlinear setting. The  
 1534 reported results are the mean and standard deviation of the metrics over 10 repetitions across different  
 1535 graph types, vanilla and model assumption violation scenarios.

1536

1537

1538

Method	$d$	SHD	SID
CAM	10	<b>8.21±5.17</b>	<b>22.66±16.98</b>
	20	<b>22.23±12.36</b>	<b>117.37±79.80</b>
	50	<b>61.92±35.50</b>	<b>696.47±463.94</b>
NOTEARS-MLP	10	12.23±3.63	44.68±10.27
	20	28.86±7.74	217.09±78.45
	50	73.40±18.41	1271.24±560.81
GraN-DAG	10	<b>10.44±5.11</b>	<b>38.19±16.78</b>
	20	30.46±8.23	214.86±75.62
	50	85.34±11.81	1453.35±575.73
DAGMA	10	12.64±4.66	46.99±17.88
	20	<b>28.28±8.83</b>	<b>211.26±85.84</b>
	50	<b>72.07±20.15</b>	<b>1254.57±583.11</b>

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552 Table 28: Summary of performances of the most competitive methods under MLP setting. The  
 1553 reported results are the mean and standard deviation of the metrics over 10 repetitions across different  
 1554 graph types, vanilla and model assumption violation scenarios.

1555

1556

1557

Method	$d$	SHD	SID
CAM	10	15.67±2.87	44.70±6.53
	20	34.88±9.23	179.14±29.38
	50	84.89±26.26	904.99±134.22
NOTEARS-MLP	10	<b>12.93±4.24</b>	<b>34.63±12.57</b>
	20	<b>26.97±8.02</b>	<b>132.86±52.22</b>
	50	<b>70.08±23.80</b>	<b>740.91±239.08</b>

1558

1559

1560

1561

1562

1563

1564

1565

1566 I TABLE RESULTS ON REAL-WORLD DATA  
1567

1568 We test the performance of 12 benchmark methods on the real-world Sachs (Sachs et al., 2005)  
1569 dataset. Sachs is a bioinformatics dataset used to study the expression levels of various proteins and  
1570 phospholipids in human cells, and it is a commonly used benchmark in the causal discovery field. We  
1571 conduct experiments based on 7466 samples. The true graph structure of the Sachs dataset contains  
1572 11 nodes and 17 edges, and it is widely accepted by the biological research community.

1573  
1574 Table 29: Results on Sachs dataset.

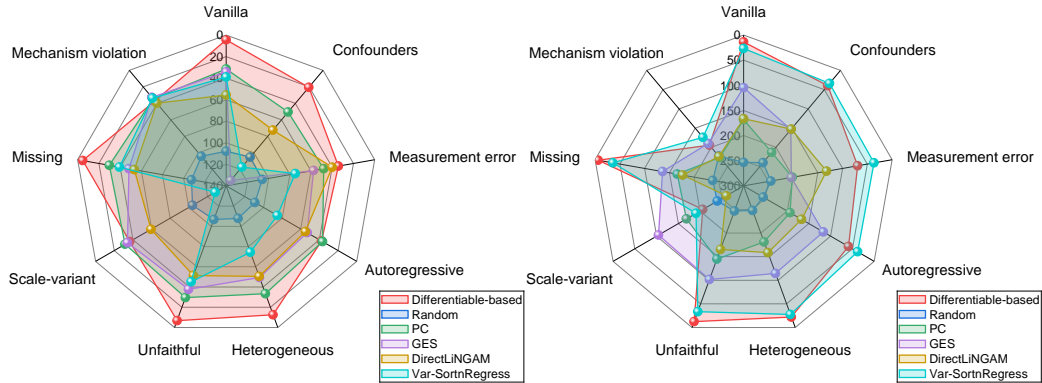
Method	SHD	SID
Random	33	56
PC	22	49
GES	30	47
DirectLiNGAM	14	50
Var-SortnRegress	19	49
$R^2$ -SortnRegress	22	51
NOTEARS	17	48
GOLEM	15	58
NoCurl	16	50
CAM	15	51
NOTEARS-MLP	14	46
GraN-DAG	15	45
DAGMA	<b>12</b>	<b>42</b>

1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590 The results in Table 29 show that, represented by DAGMA, differentiable causal discovery achieves  
1591 optimal performance on the real-world Sachs dataset. Considering that Sachs is also regarded as  
1592 a real-world heterogeneous dataset (Mooij et al., 2020), the results on both Sachs and synthetic  
1593 datasets further indicate that differentiable causal discovery performs better under model assumption  
1594 violations.

1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

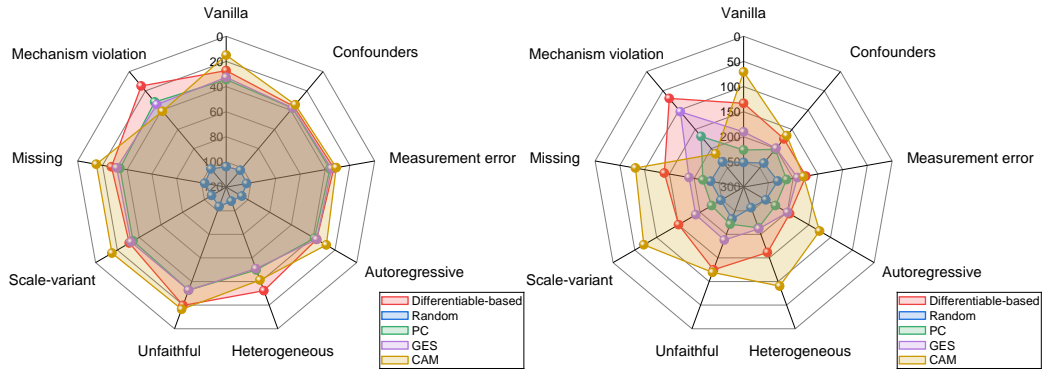
J FIGURE RESULTS ACROSS NODES, GRAPH TYPES, AND GRAPH DENSITIES

This section presents a comprehensive analysis of the figure results across varying numbers of nodes, graph types, and graph densities, in Figure 2, 3, 4, 5, 6, 7, and 8.



(a) SHD with linear ER-2 graphs of 20 nodes

(b) SID with linear ER-2 graphs of 20 nodes



(c) SHD with nonlinear ER-2 graphs of 20 nodes

(d) SID with nonlinear ER-2 graphs of 20 nodes

Figure 2: Experimental results under the linear and nonlinear ER-2 graphs of 20 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 2c and Figure 2d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

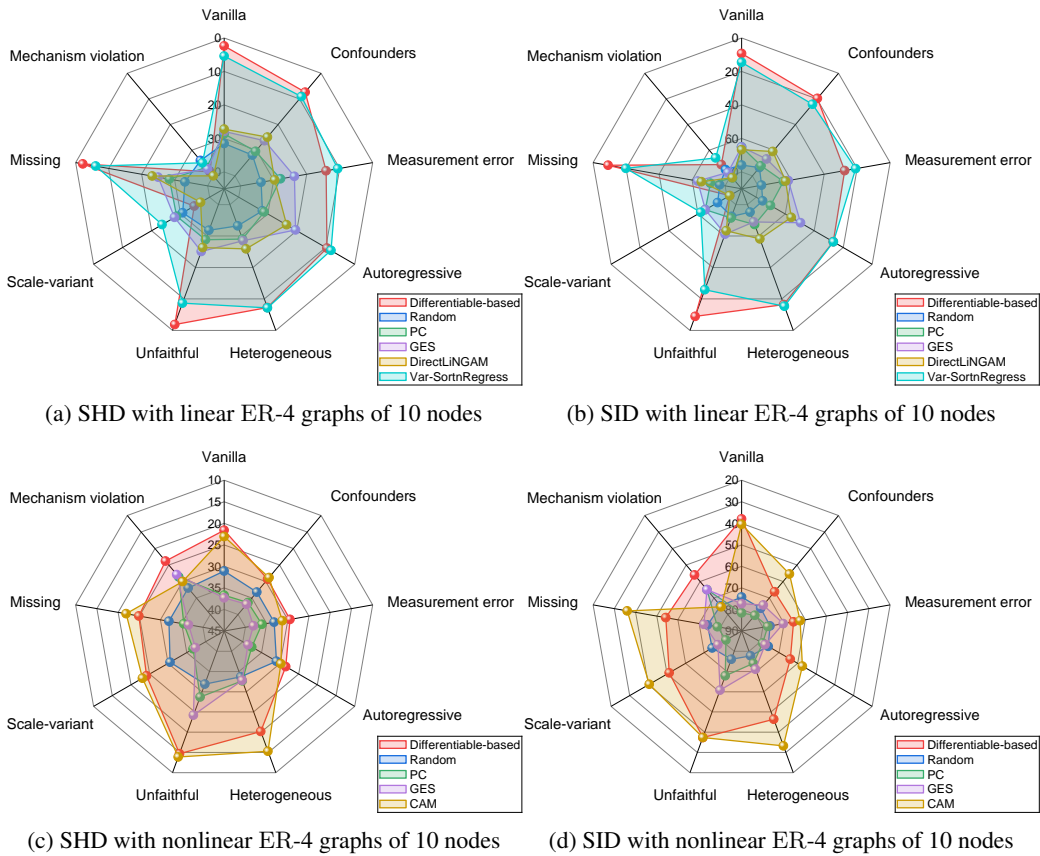


Figure 3: Experimental results under the linear and nonlinear ER-4 graphs of 10 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 3c and Figure 3d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).



1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

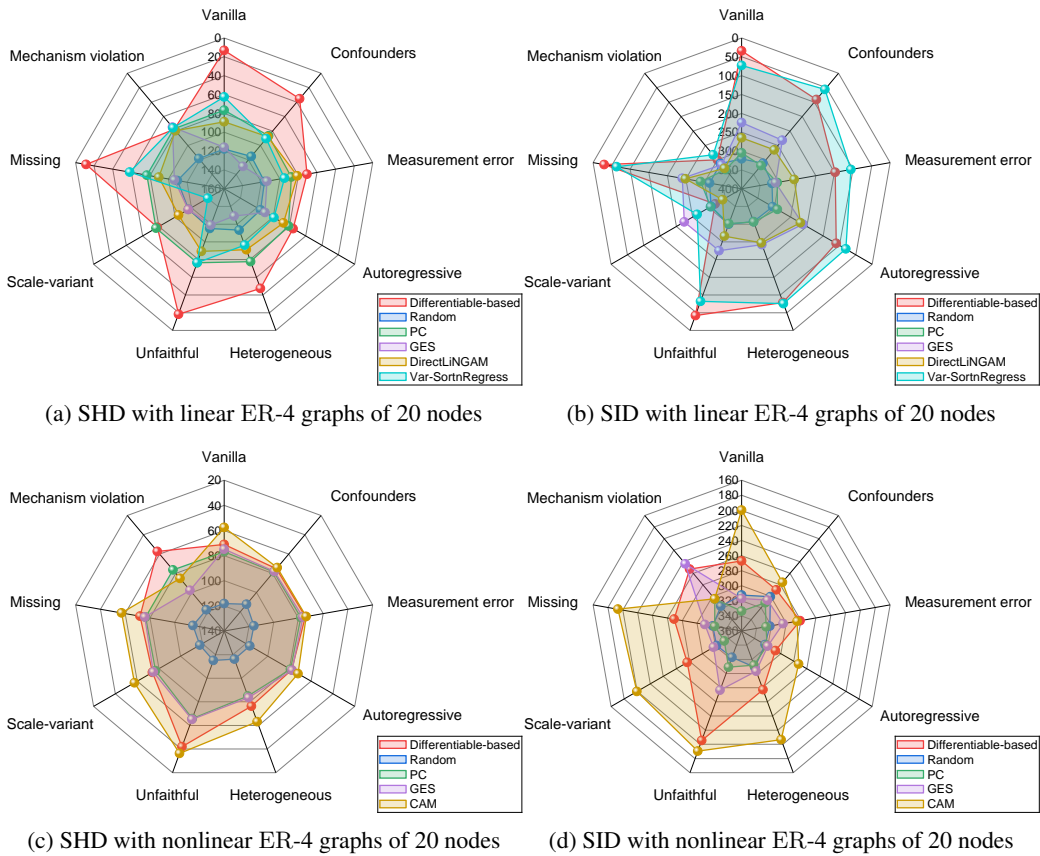


Figure 4: Experimental results under the linear and nonlinear ER-4 graphs of 20 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 4c and Figure 4d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

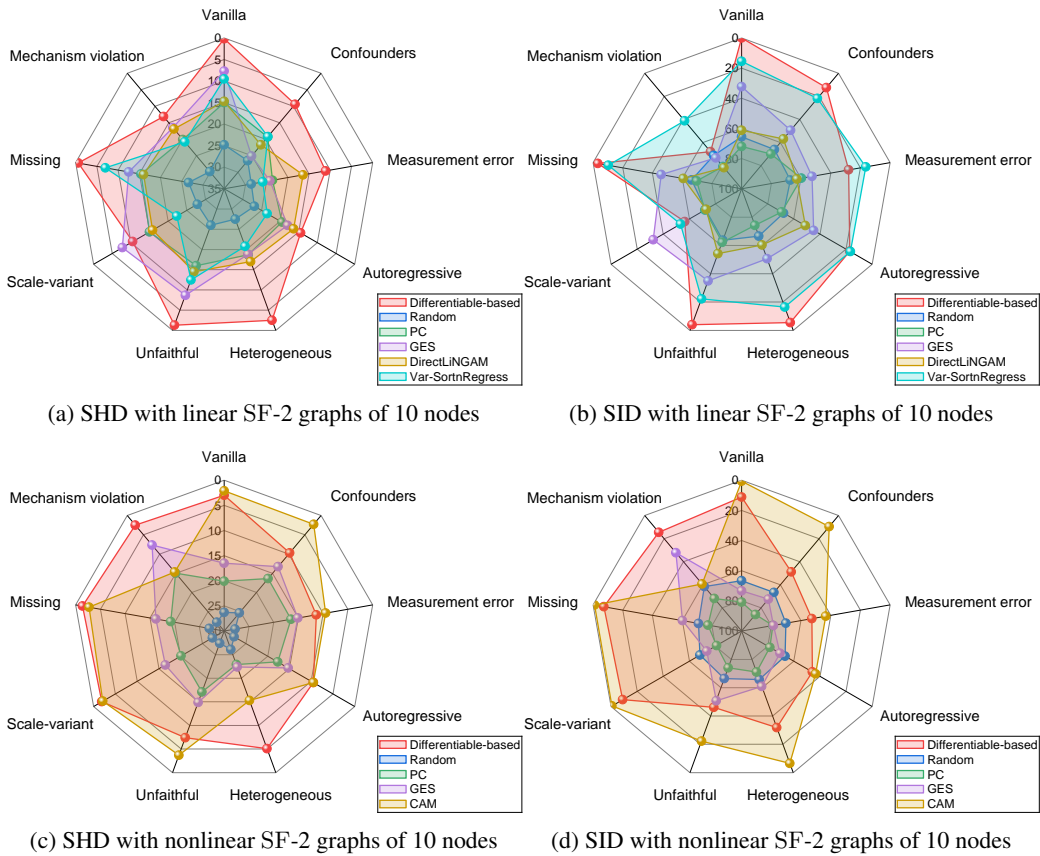


Figure 5: Experimental results under the linear and nonlinear SF-2 graphs of 10 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 5c and Figure 5d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

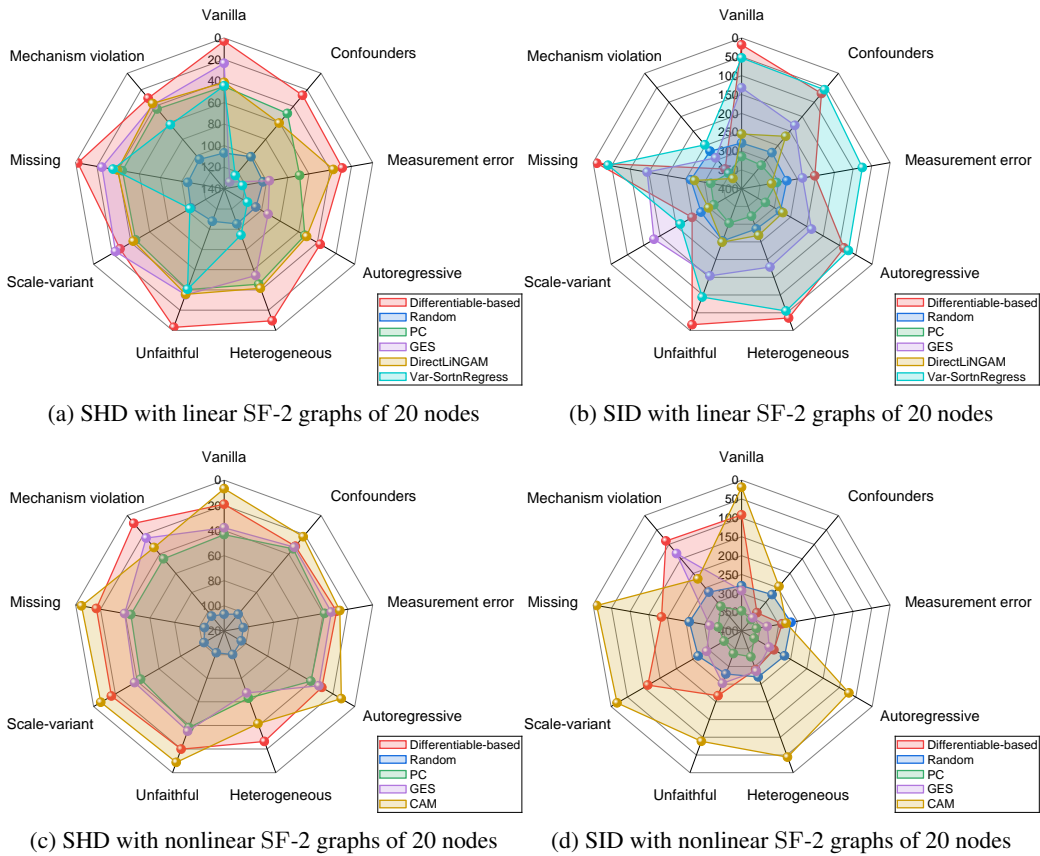


Figure 6: Experimental results under the linear and nonlinear SF-2 graphs of 20 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 6c and Figure 6d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

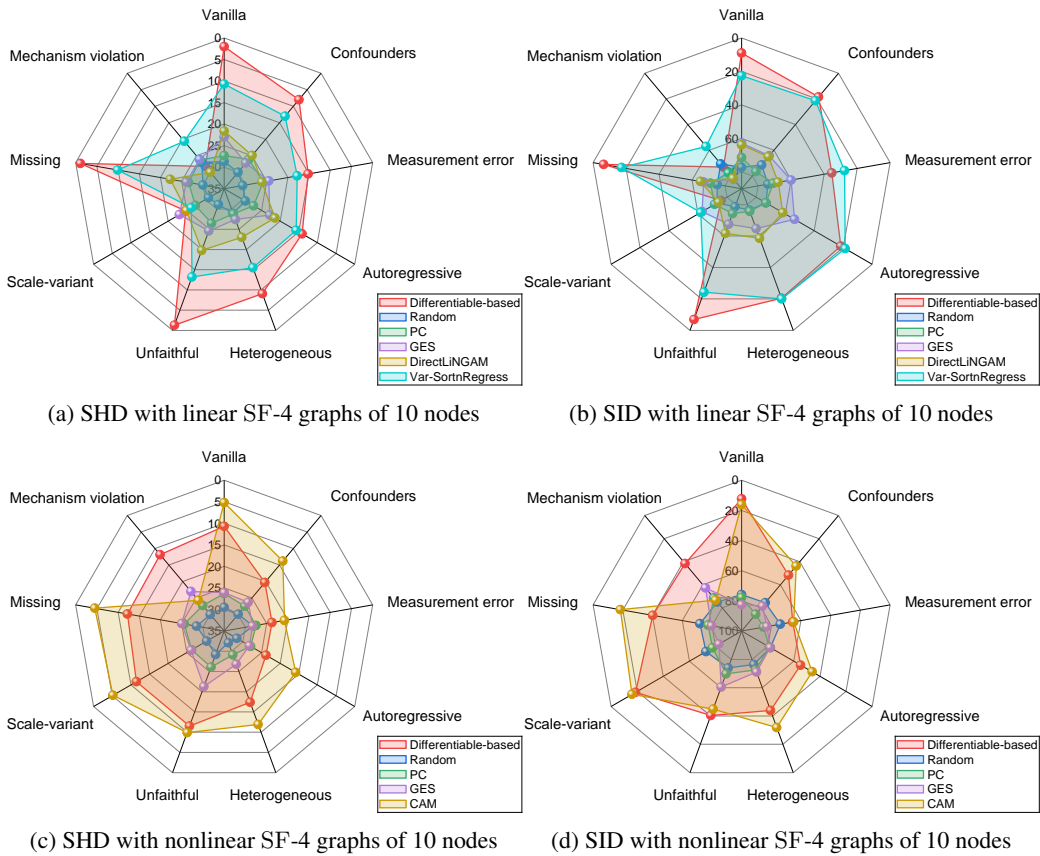


Figure 7: Experimental results under the linear and nonlinear SF-4 graphs of 10 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 7c and Figure 7d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976  
 1977  
 1978  
 1979  
 1980  
 1981  
 1982  
 1983  
 1984  
 1985  
 1986  
 1987  
 1988  
 1989  
 1990  
 1991  
 1992  
 1993  
 1994  
 1995  
 1996  
 1997

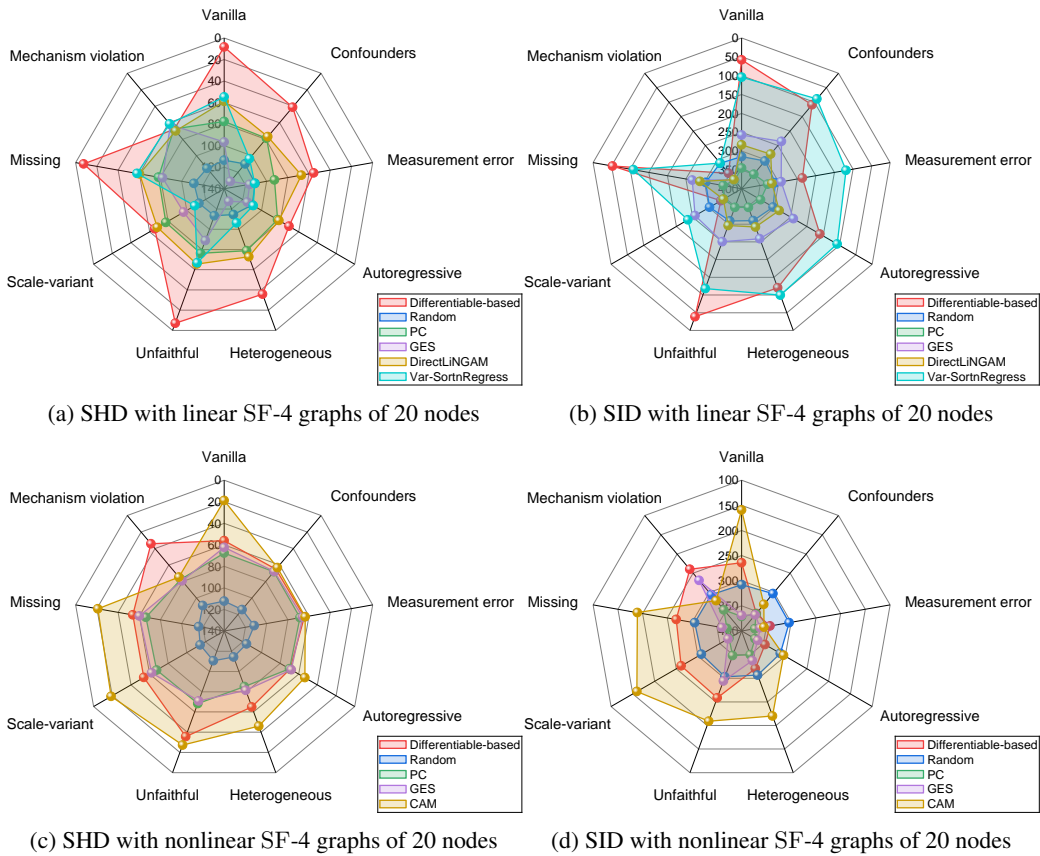


Figure 8: Experimental results under the linear and nonlinear SF-4 graphs of 20 nodes. SHD (the lower the better) and SID (the lower the better) are evaluated over 10 trials. For the differentiable causal discovery method, we present only the optimal results. As the nonlinear settings in Figure 8c and Figure 8d are more favorable to CAM, we conduct a more reasonable evaluation of CAM and differentiable causal discovery under the MLP setting (Section 4.1.1).

## K TABLE RESULTS FOR EXTREME MEASUREMENT ERROR

Table 30 presents the results in the linear setting under measurement error with  $\delta = 10$ . The results in Table 30 indicate that when  $\delta$  takes larger values, differentiable causal discovery methods fail to demonstrate robust performance.  $\delta$  is used to control the variance of  $\epsilon_i$  in (7). As  $\delta$  increases, the noise ratio correspondingly increases, which leads to the loss of robustness in differentiable methods. Tables 2.1 and 9 present the results in the linear setting under measurement error with  $\delta = 0.8$ . The results in Tables 2.1 and 9 indicate that when  $\delta = 0.8$  takes a smaller value, the noise ratio is correspondingly lower, allowing differentiable methods to demonstrate robust performance.

Table 30: Linear Setting under measurement error with  $\delta = 10$ , for ER-2 graphs of 10, 20 nodes.

10 nodes	Vanilla model		Measurement error ( $\delta = 10$ )	
	SHD↓	SID↓	SHD↓	SID↓
Random	25.6±3.1	57.9±9.5	23.1±1.9	61.2±7.5
PC	12.4±3.1	40.9±13.4	<b>19.2±2.1</b>	56.9±9.3
GES	13.8±7.8	32.0±13.6	20.2±4.5	<b>54.1±11.6</b>
DirectLiNGAM	19.6±3.3	46.1±10.6	20.0±1.1	61.2±8.2
DAGMA	<b>1.2±1.2</b>	<b>3.3±5.3</b>	20.7±1.2	58.2±7.6
20 nodes	SHD↓	SID↓	SHD↓	SID↓
Random	107.9±7.0	253.2±26.3	92.7±6.7	243.6±19.4
PC	31.6±6.5	168.5±27.6	44.5±4.6	213.8±24.5
GES	34.3±24.6	104.5±51.1	51.2±7.1	220.9±27.6
DirectLiNGAM	55.7±9.1	166.4±31.0	<b>41.8±1.8</b>	<b>210.3±24.3</b>
DAGMA	<b>5.4±3.9</b>	<b>14.2±10.3</b>	49.4±4.6	227.5±24.9

## L TABLE RESULTS ON SEMI-SYNTHETIC DATA

The semi-synthetic data is generated based on the network structure of the real-world Sachs dataset, using linear and nonlinear vanilla models to create eight datasets with model assumption violations. The results in Table 31 and 32 indicate that differentiable causal discovery methods still achieve optimal or competitive performance in scenarios other than scale variation.

Table 31: Linear Setting, for semi-synthetic data of 11 nodes.

11 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>
Random	34.0±2.0	44.1±6.9	32.7±1.4	42.9±5.3	36.2±3.1	47.2±5.8	33.7±4.2	46.8±5.7	38.4±3.7	51.6±4.3	36.8±3.7	50.3±3.9	33.5±1.8	43.2±6.5	39.1±2.8	57.3±8.2	35.9±4.8	56.7±7.4
PC	10.7±3.8	38.6±11.0	16.2±2.5	44.9±4.1	15.6±2.2	38.8±11.2	14.7±1.7	42.2±5.5	13.1±2.4	43.4±7.8	14.7±1.6	50.3±5.3	10.7±3.8	38.6±11.0	9.6±2.2	34.1±8.5	17.5±2.3	46.3±4.1
GES	9.4±2.7	28.9±6.3	28.7±7.4	33.8±12.0	17.0±3.5	29.2±14.5	14.5±3.9	22.5±14.5	15.5±3.9	27.4±9.2	13.7±3.4	35.5±8.4	<b>9.3±2.7</b>	<b>28.7±6.2</b>	7.9±2.9	26.5±8.7	19.2±3.3	35.5±11.4
DirectLINGAM	12.8±4.2	34.5±11.5	18.9±4.3	39.6±6.9	14.8±3.2	41.0±6.0	12.2±2.5	35.6±11.2	14.5±4.4	42.5±9.7	11.3±3.7	35.7±8.2	14.3±4.0	43.6±8.5	12.3±4.1	37.1±7.5	16.8±2.8	40.9±6.1
Var-SortnRegress	3.8±3.4	7.8±8.3	22.6±7.0	12.7±7.2	15.9±1.7	<b>11.3±8.3</b>	12.6±3.7	<b>8.5±9.2</b>	11.3±4.6	8.2±8.2	5.9±1.9	15.0±5.7	13.1±4.6	38.5±6.9	2.9±2.0	7.7±6.8	18.0±2.6	23.4±7.1
R <sup>2</sup> -SortnRegress	17.1±5.6	37.1±7.1	36.3±7.0	38.8±5.0	20.3±2.6	36.3±9.6	19.6±4.4	38.3±5.0	21.6±3.1	37.8±7.1	22.4±3.0	42.9±2.3	17.1±5.6	37.1±7.1	18.1±6.5	40.4±7.6	22.4±2.0	39.5±5.7
NOTEARS	0.5±0.7	5.1±6.3	9.3±1.6	29.0±6.2	<b>8.9±2.4</b>	<b>18.3±7.5</b>	11.5±4.4	11.3±8.8	3.4±1.9	10.5±7.1	0.4±0.9	3.2±6.5	12.1±3.6	44.3±9.2	1.1±1.6	5.2±7.5	14.6±2.4	31.7±5.0
GOLEM	0.7±0.3	6.8±5.4	9.3±2.6	26.7±3.1	11.3±4.9	25.7±5.9	15.0±2.8	37.7±5.4	5.7±5.2	<b>2.7±3.8</b>	0.7±0.9	2.3±3.3	13.0±1.4	46.3±3.4	1.5±0.8	7.2±4.3	17.0±0.8	52.3±0.5
NoCurl	0.3±0.6	2.8±5.9	12.0±3.9	<b>9.8±7.2</b>	9.4±2.3	<b>18.3±8.0</b>	11.7±4.3	<b>9.2±7.5</b>	3.1±2.3	6.1±5.1	0.5±0.9	4.6±7.1	12.2±2.4	<b>43.0±8.4</b>	<b>0.5±0.7</b>	4.1±5.4	22.8±3.5	<b>22.8±6.1</b>
DAGMA	<b>0.2±0.4</b>	<b>2.8±5.6</b>	<b>8.9±2.8</b>	25.9±10.3	<b>8.9±2.4</b>	19.1±8.3	<b>10.9±3.9</b>	12.9±10.7	<b>2.9±1.9</b>	6.3±6.0	<b>0.1±0.3</b>	<b>1.4±4.2</b>	<b>11.5±2.8</b>	43.7±10.1	0.6±1.0	<b>2.4±3.7</b>	<b>13.9±1.9</b>	28.0±6.4

Table 32: MLP Setting, for semi-synthetic data of 11 nodes.

11 nodes	Vanilla model		Latent confounders		Measurement error		Autoregressive		Heterogeneous		Unfaithful		Scale-variant		Missing		Mechanism violation	
	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>	SHD <sub>i</sub>	SID <sub>i</sub>
Random	32.9±3.2	44.0±7.9	30.8±2.7	41.9±5.8	37.5±6.4	58.9±8.2	35.2±2.9	47.3±4.8	31.6±2.9	46.8±4.7	33.9±1.5	54.7±5.6	36.2±2.8	57.3±4.6	38.7±6.3	51.7±5.4	31.4±2.3	45.9±5.1
PC	17.5±2.3	46.3±4.1	18.8±3.9	48.6±4.7	18.4±1.9	37.1±10.7	17.1±2.1	44.4±4.8	20.9±2.9	44.3±6.1	18.2±2.6	45.7±5.4	17.5±2.3	46.3±4.1	16.8±1.5	44.6±4.8	10.7±3.8	38.6±11.0
GES	19.2±3.3	35.5±11.4	23.1±5.6	42.9±11.2	19.6±3.2	45.2±10.3	17.6±4.5	34.0±11.6	27.5±4.0	35.2±5.8	21.1±3.2	40.9±9.5	19.5±3.1	35.5±11.7	19.6±2.0	37.0±10.4	9.4±2.7	28.9±6.3
CAM	9.4±3.4	13.5±10.7	15.7±4.1	39.1±7.2	21.9±5.5	42.3±11.4	13.0±3.8	35.7±12.6	18.1±3.5	28.6±12.4	12.3±2.7	20.1±8.3	<b>9.4±3.4</b>	<b>13.5±10.7</b>	9.3±3.6	12.6±9.8	13.4±1.7	43.3±6.4
NOTEARS-MLP	<b>6.8±2.9</b>	<b>10.1±7.5</b>	10.4±1.6	<b>38.2±7.3</b>	15.4±1.2	46.5±3.9	15.9±4.1	39.9±5.2	13.9±2.9	<b>20.7±9.1</b>	10.0±2.9	<b>17.3±8.3</b>	16.5±2.3	46.5±9.8	<b>7.3±2.4</b>	<b>7.2±5.8</b>	5.2±1.7	<b>22.2±5.7</b>
Grn-DAG	9.1±2.6	33.8±10.0	<b>10.0±2.6</b>	40.9±8.8	<b>13.0±2.1</b>	42.0±3.9	<b>10.5±2.5</b>	<b>33.3±7.5</b>	<b>9.1±1.8</b>	32.6±7.7	<b>10.0±1.5</b>	28.1±4.3	13.1±1.5	47.1±6.7	9.8±2.7	35.7±10.1	10.6±2.5	40.0±7.3
DAGMA	8.7±2.7	11.1±5.1	11.9±1.7	41.0±5.6	15.8±2.6	<b>34.7±3.8</b>	16.0±1.5	40.5±8.9	16.2±2.7	34.7±8.4	13.2±2.4	21.5±8.5	16.7±2.2	43.5±4.8	9.0±2.8	12.1±5.9	<b>5.0±1.4</b>	22.6±4.2