
ANALYTIC DAG CONSTRAINTS FOR DIFFERENTIABLE DAG LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Recovering underlying Directed Acyclic Graph (DAG) structures from observational data presents a formidable challenge due to the combinatorial nature of the DAG-constrained optimization problem. Recently, researchers have identified gradient vanishing as one of the primary obstacles in differentiable DAG learning and have proposed several DAG constraints to mitigate this issue. By developing the necessary theory to establish a connection between analytic functions and DAG constraints, we demonstrate that analytic functions from the set $\{f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i \mid \forall i > 0, c_i > 0; r = \lim_{i \rightarrow \infty} c_i / c_{i+1} > 0\}$ can be employed to formulate effective DAG constraints. Furthermore, we establish that this set of functions is closed under several functional operators, including differentiation, summation, and multiplication. Consequently, these operators can be leveraged to create novel DAG constraints based on existing ones. Using these properties, we designed a series of DAG constraints and developed an efficient algorithm to evaluate these DAG constraints. Experiments conducted in various settings demonstrate that our DAG constraints outperform previous state-of-the-art approaches.

1 INTRODUCTION

DAG learning aims to recover Directed Acyclic Graphs (DAGs) from observational data, which is a core problem in many fields, including bioinformatics (Sachs et al., 2005; Zhang et al., 2013), machine learning (Koller and Friedman, 2009), and causal inference (Spirtes et al., 2000). Under certain assumptions (Pearl, 2000; Spirtes et al., 2000), the recovered DAGs [could be interpreted causally](#) (Koller and Friedman, 2009) and hold causal interpretations.

There are two main categories of DAG learning approaches: constraint-based and score-based methods. Most constraint-based approaches, *e.g.*, PC (Spirtes and Glymour, 1991), FCI (Spirtes et al., 1995; Colombo et al., 2012), rely on conditional independence tests, which typically necessitate a large sample size (Shah and Peters, 2020; Vowels et al., 2021). The score-based approaches, including exact methods based on dynamic programming (Koivisto and Sood, 2004; Singh and Moore, 2005; Silander and Myllymäki, 2006), A* search (Yuan et al., 2011; Yuan and Malone, 2013), and integer programming (Cussens, 2011), as well as greedy methods like GES (Chickering, 2002), model the validity of a graph according to some score function and are often formulated and solved as discrete optimization problems. A key challenge for score-based methods is the super-exponential combinatorial search space of DAGs w.r.t number of nodes (Chickering, 1996; Chickering et al., 2004).

Recently, Zheng et al. (2018) developed a continuous DAG learning approach using Lagrange Multiplier methods and a differentiable DAG constraint based on the trace of the matrix exponential of the weighted adjacency matrix. The resulting method, named NOTEARS, demonstrated superior performance in estimating linear DAGs with equal noise variances. Very recently, Zhang et al. (2022) and Bello et al. (2022) suggest that one main issue for NOTEARS and its derivatives, such as Yu et al. (2019), is gradient vanishing for linear DAG models with equal variance. They have thus proposed new continuous DAG constraints by based on geometric series of matrices as well as log-determinant of matrices.

In fact, many of the proposed Directed Acyclic Graph (DAG) constraints can be unified, as demonstrated in Wei et al. (2020). [Wei et al. \(2020\)](#) reveals that, for a $d \times d$ adjacency matrix, an order- d

polynomial of matrices is necessary and sufficient to enforce the DAG property. However, from a computational standpoint, computing general matrix polynomials can be challenging. Considering the fact that infinite-order polynomials that converge, i.e., power series, can give rise to analytic functions that are often simpler to evaluate than general polynomials, it prompts the question of whether analytic functions could be utilized in constructing DAG constraints. Furthermore, it raises the possibility of employing techniques commonly used for analyzing analytic functions in the investigation of continuous DAG constraints.

The answer is yes. We demonstrate that any analytic function within the class of functions denoted as $\mathcal{F} = \{f|f(x) = c_0 + \sum_{i=0}^{\infty} c_i x^i; c_i > 0, \forall i > 0; \lim_{i \rightarrow \infty} c_i/c_{i+1} > 0\}$ can be utilized to formulate Directed Acyclic Graph (DAG) constraints. In fact, the DAG constraints introduced in Zheng et al. (2018), Zhang et al. (2022), and Bello et al. (2022) can all be interpreted as being based on analytic functions from \mathcal{F} . Furthermore, we establish that the function class \mathcal{F} remains closed under various function operators, including differentiation, function addition, and function multiplication. Leveraging this insight, we can construct novel DAG constraints based on pre-existing ones. Additionally, we can analyze the performance of these derived DAG constraints using techniques rooted in analytic functions.

2 PRELIMINARIES

DAG and Linear SEM Given a directed acyclic graph (DAG) \mathcal{G} defined over random vector $\mathbf{x} = [x_1, x_2, \dots, x_d]^\top$, the corresponding distribution $P(\mathbf{x})$ is assumed to satisfy the Markov assumption (Spirtes et al., 2000; Pearl, 2000). We consider \mathbf{x} to follow a linear Structural Equation Model (SEM):

$$\mathbf{x} = \mathbf{B}^\top \mathbf{x} + \mathbf{e}. \quad (1)$$

Here, $\mathbf{B} \in \mathbb{R}^{d \times d}$ represents the weighted adjacency matrix that characterizes the DAG \mathcal{G} , and $\mathbf{e} = [e_1, e_2, \dots, e_d]^\top$ represents the exogenous noise vector, comprising d independent random variables. To simplify notation, we use $\mathcal{G}(\mathbf{B})$ to denote the graph induced by the weighted adjacency matrix \mathbf{B} , and we interchangeably use the terms ‘random variables’ and ‘vertices’ or ‘nodes’.

We aim to estimate the DAG \mathcal{G} from n i.i.d. observational examples of \mathbf{x} , denoted by $\mathbf{X} \in \mathbb{R}^{n \times d}$. Generally, the DAG \mathcal{G} can be identified only up to its Markov equivalence class under the faithfulness (Spirtes et al., 2000) or the sparsest Markov representation assumption (Raskutti and Uhler, 2018). It has been demonstrated that for linear SEMs with homoscedastic errors, where the noise terms are specified up to a constant (Loh and Bühlmann, 2013), and for linear non-Gaussian SEMs, where no more than one of the noise terms is Gaussian (Shimizu et al., 2006), the true DAG can be fully identified. In our study, we specifically focus on linear SEMs with equal noise variances (Peters and Bühlmann, 2013), where the scale of the data may be either known or unknown. When the scale is known, it is possible to fully recover the DAG. However, in the case of an unknown scale, the DAG may only be identified up to its Markov equivalence class.

Continuous DAG learning In recent years, a series of continuous Directed Acyclic Graph (DAG) learning algorithms Bello et al. (2022); Ng et al. (2020); Zhang et al. (2022); Yu et al. (2021; 2019); Zheng et al. (2018) has been introduced, demonstrating superior performance when applied to linear Structural Equation Models (SEMs) with equal noise variances and known data scale. These methods can be expressed as follows:

$$\operatorname{argmin}_{\mathbf{B}} S(\mathbf{B}, \mathbf{X}), \text{ s.t. } h(\mathbf{B}) = 0. \quad (2)$$

Here, S is a scoring function, which can take the form of mean square error (Zheng et al., 2018) or negative log-likelihood (Ng et al., 2020). The function h is continuous and equal to 0 if and only if the weighted adjacency matrix \mathbf{B} defines a valid DAG. Previous approaches have employed various techniques, such as matrix exponential (Zheng et al., 2018), log-determinants (Bello et al., 2022), and polynomials (Zhang et al., 2022), to construct the function h . However, these methods are known to perform poorly when applied to normalized data since they rely on scale information across variables for complete DAG recovery (Reisach et al., 2021).

3 ANALYTIC DAG CONSTRAINTS

In this section, we demonstrate that the diverse set of continuous DAG constraints proposed in previous work can be unified through the use of analytic functions. We will begin by offering a brief introduction to analytic functions and then illustrate how they can be employed to establish DAG constraints.

3.1 ANALYTIC FUNCTIONS AS DAG CONSTRAINTS

In mathematics, a power series

$$f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i, \quad (3)$$

which converges for $|x| < r = \lim_{i \rightarrow \infty} |c_i/c_{i+1}|$, defines an analytic function f on the open interval $(-r, r)$, and r is known as the convergence radius. When we replace x with a square matrix \mathbf{A} , we obtain an analytic function f of a matrix as follows:

$$f(\mathbf{A}) = c_0 \mathbf{I} + \sum_{i=1}^{\infty} c_i \mathbf{A}^i, \quad (4)$$

where \mathbf{I} is the identity matrix. Equation (4) would converge if the largest absolute value of eigenvalues of \mathbf{A} , known as the spectral radius and denoted by $\rho(\mathbf{A})$, is smaller than r .

We are particularly interested in the following specific class of analytic functions

$$\mathcal{F} = \{f | f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i; \forall i > 0, c_i > 0; \lim_{i \rightarrow \infty} c_i/c_{i+1} > 0\}, \quad (5)$$

as any analytic function belongs to \mathcal{F} can be applied to construct a continuous DAG constraint.

Proposition 1. *Let $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with $\rho(\tilde{\mathbf{B}}) < r$ be the weighted adjacency matrix of a directed graph \mathcal{G} , and let f be an analytic function in the form of (4), where we further assume $\forall i > 0$ we have $c_i > 0$, then \mathcal{G} is acyclic if and only if*

$$\text{tr} [f(\tilde{\mathbf{B}})] = c_0 d. \quad (6)$$

An interesting property the DAG constraint (6) is that its gradients can also be represented as transpose of an analytic function as follows, which allows us to use analytic functions as the gradients of DAG constraints.

Proposition 2. *There exists some real number r , where for all $\{\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\tilde{\mathbf{B}}) < r\}$, the derivative of $\text{tr} [f(\tilde{\mathbf{B}})]$ w.r.t. $\tilde{\mathbf{B}}$ is*

$$\nabla_{\tilde{\mathbf{B}}} \text{tr} [f(\tilde{\mathbf{B}})] = [\nabla_x f(x)|_{x=\tilde{\mathbf{B}}}]^{\top}. \quad (7)$$

It is notable that for a $d \times d$ weighted adjacency matrix $\tilde{\mathbf{B}}$, an order- d polynomial of $\tilde{\mathbf{B}}$ is sufficient and necessary to enforce DAGness (Wei et al., 2020; Ng et al., 2022). Meanwhile, evaluating matrix polynomials efficiently is highly nontrivial (Higham, 2008). For matrix analytic functions such as exponentials or logarithms, however, efficient algorithms exist (Higham, 2008).

The connection between matrix analytic functions and real analytic functions means that various properties of the matrix function can be obtained from a simple real-valued function. To pursue DAG constraints with better computational efficiency, we seek an analytic function whose derivative can be represented by itself to reduce the computation of different analytic functions. If a function has such property, various intermediate results can be saved for future computation of gradients. The exponential function $\exp(x)$ with $\partial \exp(x)/\partial x = \exp(x)$, is a natural contender, and this leads to the well-known exponential-based DAG constraints (Zheng et al., 2018)

$$\text{Constraints: } \text{tr} [\exp(\tilde{\mathbf{B}})] = \sum_{i=0}^{\infty} \tilde{\mathbf{B}}^i / i! = d, \quad \text{Gradient: } \nabla_{\tilde{\mathbf{B}}} \exp(\tilde{\mathbf{B}}) = \exp(\tilde{\mathbf{B}})^{\top}, \quad (8)$$

162 which will converge for any $\tilde{\mathbf{B}}$.

163
164 Recently Bello et al. (2022) and Zhang et al. (2022) have suggested that exponential-based DAG
165 constraints suffers from gradient vanishing. One cause of gradient vanishing arises from the
166 small coefficients of high order terms. The convergence radius for the exponential is ∞ , that
167 is $\lim_{i \rightarrow \infty} |c_i/c_{i+1}| = \lim_{i \rightarrow \infty} |(i+1)!/i!| = \infty$, which suggests that, compared to the lower order terms,
168 the higher order terms contribute almost nothing in the DAG constraints, which indicates that it would
169 not be efficient to prohibit possible long loops in candidate adjacency matrices.

170 Due to the fact that the adjacency matrix of a DAG must form a nilpotent matrix, whose spectral
171 radius are acutally 0, naturally the spectral radius of candidate adjacency matrices would be close to
172 0. As a result, we do not need a function with infinite convergence radius. Instead, we can use an
173 analytic function with finite convergence radius $r = \lim_{i \rightarrow \infty} |c_i/c_{i+1}| < \infty$. Thus by using a sequence
174 c_i with geometric progression $c_i = 1/s^{i-1}$ or harmonic-geometric progression $c_i = 1/(is^{i-1})$ we
175 can obtain two analytic functions,

$$176 f_{inv}^s(x) = (s-x)^{-1} = \sum_{i=0}^{\infty} x^i/s^{i-1}, \quad f_{log}^s(x) = -s \log(s-x) = \sum_{i=1}^{\infty} \frac{x^i}{is^{i-1}} - s \log s. \quad (9)$$

177 Then by our Proposition 1 and Proposition 2, two dag constraints can be obtained as follows:

$$181 \text{Constraints: } \text{tr} f_{inv}^s(\tilde{\mathbf{B}}) = d, \quad \text{Gradient: } \nabla_{\tilde{\mathbf{B}}} \text{tr} f_{inv}^s(\tilde{\mathbf{B}}) = [f_{inv}^s(\tilde{\mathbf{B}})^2]^\top, \quad (10a)$$

$$182 \text{Constraints: } \text{tr} f_{log}^s(\tilde{\mathbf{B}}) = 0, \quad \text{Gradient: } \nabla_{\tilde{\mathbf{B}}} \text{tr} f_{log}^s(\tilde{\mathbf{B}}) = [f_{inv}^s(\tilde{\mathbf{B}})]^\top, \quad (10b)$$

183 where a truncated version of f_{inv}^s is applied in Zhang et al. (2022), and the f_{log}^s based constraints are
184 equivalent to those in Bello et al. (2022). One key difference between Zhang et al. (2022); Bello et al.
185 (2022) and the exponential-based DAG constraints (Zheng et al., 2018) is their finite convergence
186 radius, which requires an additional constraints $\rho(\tilde{\mathbf{B}}) < s$. Meanwhile, the adjacency matrix of a
187 DAG must be nilpotent, and thus its spectral radius must be 0. In this case, such additional constraints
188 would not affect the feasible set.

190 3.2 CONSTRUCTING DAG CONSTRAINTS BY FUNCTIONAL OPERATOR

191 One can easily observe a coincidence between f_{log} and f_{inv} as follows,

$$192 \frac{\partial f_{log}^s(x)}{\partial x} = f_{inv}^s(x), \quad f_{log}^s(x) = \int f_{inv}^s(t) dt + C, \quad (11)$$

193 which suggests that it may be possible to derive a group of DAG constraints from an analytic function
194 by applying integration or differentiation. This is because derivatives of any order of an analytic
195 function is also analytic. More formally, if a function is analytic at some point x_0 , then its n^{th}
196 derivative for any integer n exists and is also analytic at x_0 . Thus we can derive DAG constraints
197 from any $f \in \mathcal{F}$ as follows.

201 **Proposition 3.** Let $f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i \in \mathcal{F}$ be analytic on $(-r, r)$, and let n be arbitrary integer
202 larger than 1, then $\tilde{\mathbf{B}} \in \mathbb{R}_{>0}^{d \times d}$ with spectral radius $\rho(\tilde{\mathbf{B}}) \leq r$ forms a DAG if and only if

$$203 \text{tr} \left[\frac{\partial^n f(x)}{\partial x^n} \Big|_{x=\tilde{\mathbf{B}}} \right] = n!c_n. \quad (12)$$

204 The above proposition suggests that the differential operator can be applied to an analytic function to
205 form a new DAG constraints. Besides the differential operator, the addition and multiplication of
206 analytic functions can also be applied to generate new DAG constraints. That is

207 **Proposition 4.** Let $f_1(x) = c_0^1 + \sum_{i=1}^{\infty} c_i^1 x^i \in \mathcal{F}$, and $f_2(x) = c_0^2 + \sum_{i=1}^{\infty} c_i^2 x^i \in \mathcal{F}$. Then for an ad-
208 jacency matrix $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with spectral radius $\rho(\tilde{\mathbf{B}}) \leq \min(\lim_{i \rightarrow \infty} c_i^1/c_{i+1}^1, \lim_{i \rightarrow \infty} c_i^2/c_{i+1}^2)$,
209 the following three statements are equivalent:

- 210 1. $\tilde{\mathbf{B}}$ forms a DAG;
- 211 2. $\text{tr}[f_1(\tilde{\mathbf{B}}) + f_2(\tilde{\mathbf{B}})] = (c_0^1 + c_0^2)d$;

$$3. \text{tr}[f_1(\tilde{\mathbf{B}})f_2(\tilde{\mathbf{B}})] = c_0^1 c_0^2 d.$$

Particularly for $f_{\log}^s(x)$ and $f_{\text{inv}}^s(x)$, due to the specific property of $f_{\text{inv}}^s(x)$, we have

$$\frac{\partial^{n+1} f_{\log}^s(x)}{\partial x^{n+1}} = \frac{\partial^n f_{\text{inv}}^s(x)}{\partial x^n} \propto (s-x)^{-(n+1)} = [f_{\text{inv}}^s(x)]^{n+1}, \quad (13)$$

which implies that the n^{th} derivative of function $1/(s-x)$ is propositional to the the order- $(n+1)$ power of $1/(s-x)$. Using this property, the value of $(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-1}$ can be cached and then used to generate a series of DAG constraints as well as their gradients. Similarly, the value of matrix exponential $\exp(\tilde{\mathbf{B}}/s)$ can also be cached during the evaluation of DAG constraints to accelerate the computation. Furthermore, the gradients of the DAG constraints will also increase as n increases.

Proposition 5. *Let n be any positive integer, the adjacency matrix $\tilde{\mathbf{B}} \in \{\hat{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\hat{\mathbf{B}}) < s\}$ forms a DAG if and only if*

$$\text{tr}[(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n}] = d.$$

Furthermore, the gradients of the DAG constraints satisfies that $\forall \tilde{\mathbf{B}} \in \{\hat{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\hat{\mathbf{B}}) < s\}$

$$\|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n}\| \leq \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n-k}\|,$$

where k is an arbitrary positive integer, and $\|\cdot\|$ denote an arbitrary matrix norm induced by vector p -norm.

Gradient Vanishing and Numeric Stability For the series of DAG constraints constructed from Equation (13), as gradient vanishing is one of the main challenges for differentiable DAG learning, according to Proposition 5 we may prefer larger n to achieve better performance in practice. Furthermore, choosing a smaller s may also help to amplify the gradient of DAG constraints. Therefore, Bello et al. (2022) applied an annealing strategy on s to improve performance, while Zhang et al. (2022) used a fixed $s = 1.0$ in their implementation. However, in practice, especially when incorporating the DAG constraints with first-order optimizers, the spectral radius of the candidate $\tilde{\mathbf{B}}$ can often be larger than s . Bello et al. (2022) applied a simple heuristics to search for the proper s , while Zhang et al. (2022) truncated the power series to avoid numerical issues in higher-order terms. However, in practice, we observed that Zhang et al. (2022) encountered some numerical issues for large graphs, and the simple heuristics used by Bello et al. (2022) may result in a sacrifice in performance. Based on our analysis, if $\tilde{\mathbf{B}}$ goes out, it can be verified by checking if the power series $\sum_{i=0}^{\infty} (\tilde{\mathbf{B}}/s)^i$ converges to $(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-1}$, and s can be chosen based on the spectral radius of $\tilde{\mathbf{B}}$.

Algorithm 1 Efficient Evaluation of Gradients

Input: $\tilde{\mathbf{B}}, s, d, \epsilon > 0, \xi > 0$
Output: $\nabla_{\tilde{\mathbf{B}}} \text{tr} f_{\log}^s(\tilde{\mathbf{B}})$ or $\nabla_{\tilde{\mathbf{B}}} \text{tr} [f_{\text{inv}}^s(\tilde{\mathbf{B}})]^n$

- 1: $\mathbf{D} \leftarrow \mathbf{I} + \tilde{\mathbf{B}}/s, \mathbf{W} \leftarrow \tilde{\mathbf{B}}/s$
- 2: $k = 1$
- 3: **while** $\|\mathbf{D}(\mathbf{I} - \tilde{\mathbf{B}}) - \mathbf{I}\| > \epsilon$ and $k < 2d$ **do**
- 4: $\mathbf{W} \leftarrow \mathbf{W} \times \mathbf{W}$
- 5: $\mathbf{D} \leftarrow \mathbf{D} \times (\mathbf{W} + \mathbf{I})$
- 6: $k \leftarrow 2k$
- 7: **end while**
- 8: **if** $\|\mathbf{D}(\mathbf{I} - \tilde{\mathbf{B}}) - \mathbf{I}\| > \epsilon$ **then**
- 9: $s \leftarrow \rho(\tilde{\mathbf{B}}) + \xi$, goto line 1
- 10: **else**
- 11: For f_{\log}^s return \mathbf{D}^\top / s
- 12: For $[f_{\text{inv}}^s(\tilde{\mathbf{B}})]^n$ return $n[\mathbf{D}^\top / s]^{n+1}$
- 13: **end if**

Algorithm 2 Path following algorithm

Input: $\mathbf{X} \in \mathbb{R}^{n \times d}; S; f \in \mathcal{F}; \lambda_1; \mu_0;$
 $\alpha \in (0, 1); T_{\text{outer}}; T_{\text{inner}}; \gamma > 0$
Output: Estimated \mathbf{B}

- 1: $i \leftarrow 0, \mu \leftarrow \mu_0, \mathbf{B}_0 = \mathbf{0}$
- 2: **for** $i = 0; i < T_{\text{outer}}; i++$ **do**
- 3: $\mathbf{B}_{i+1} \leftarrow \mathbf{B}_i$ ▷ Optimize over
- 4: **for** $j = 0; j < T_{\text{inner}}; j++$ **do**
- 5: $\tilde{\mathbf{B}} \leftarrow \mathbf{B}_{i+1} \odot \mathbf{B}_{i+1}$
- 6: $\mathbf{B}_{i+1} \leftarrow \mathbf{B}_{i+1} - \gamma \mu [\nabla_{\tilde{\mathbf{B}}} S(\mathbf{B}, \mathbf{X}) + \lambda_1 \text{sign}(\mathbf{B})] - \gamma \nabla_{\tilde{\mathbf{B}}} f(\tilde{\mathbf{B}}) \odot \mathbf{B}_{i+1}$
- 7: **end for**
- 8: $\mu \leftarrow \mu \times \alpha$
- 9: $\mathbf{B} \leftarrow \mathbf{B}_{i+1}$
- 10: **end for**
- 11: **Return** $\hat{\mathbf{B}}$

Efficiently Computation The specific structure of the power series $\sum_{i=0}^{\infty} (\tilde{\mathbf{B}}/s)^i$ allows for fast evaluation. Let

$$\mathbf{L}_t = \sum_{i=0}^t (\tilde{\mathbf{B}}/s)^i, \quad (14)$$

then it is evident that

$$\mathbf{L}_{2t} = \mathbf{L}_t + (\tilde{\mathbf{B}}/s)^t \mathbf{L}_t, \quad (15)$$

which indicates that the term \mathbf{L}_t can be obtained with $\mathcal{O}(\log t)$ time complexity. Furthermore, using Equation (13), the gradient of $\text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n}$ can also be easily derived from \mathbf{L}_{∞} . Along with the strategy for searching s , we can use Algorithm 1 to efficiently compute the DAG constraints.

3.3 OVERALL OPTIMIZATION FRAMEWORK

The DAG constraints above are applicable only to positive adjacency matrices, so we use the Hadamard product to map a real adjacency matrix to a positive one. Thus Equation (2) becomes:

$$\underset{\mathbf{B}}{\text{argmin}} S(\mathbf{B}, \mathbf{X}), \quad \text{s.t. } \text{tr} f(\mathbf{B} \odot \mathbf{B}) = c_0 d, \rho(\mathbf{B} \odot \mathbf{B}) < r, \quad (16)$$

where the analytic function $f(x) = c_0 + \sum_{i=1}^{\infty} x^i \in \mathcal{F}$, and \odot denotes the Hadamard product.

In our work, we choose to use the path-following approach with an ℓ_1 regularizer, as in Bello et al. (2022). This is because in the Lagrange approaches applied in Zhang et al. (2022); Yu et al. (2021); Zheng et al. (2018); Yu et al. (2019), the Lagrangian multiplier must be set to very large value to enforce DAGness, which may result in numerical instability. In the path-following approach, instead of using large Lagrangian multipliers, a small coefficients are added to the score function S as follows¹

$$\underset{\mathbf{B}}{\text{argmin}} \mu[S(\mathbf{B}, \mathbf{X}) + \lambda_1 \|\mathbf{B}\|_1] + \text{tr} f(\mathbf{B} \odot \mathbf{B}), \quad \text{s.t. } \rho(\mathbf{B} \odot \mathbf{B}) < r, \quad (17)$$

where λ_1 is the user-specified weight for the ℓ_1 regularizer. For the additional constraints $\rho(\mathbf{B} \odot \mathbf{B}) < r$, with properly chosen initial value and step-length, it can usually be satisfied. Also it is notable that $\|\mathbf{B}\|_1 < r$ is a sufficient condition for $\rho(\mathbf{B} \odot \mathbf{B}) < r$, and thus the sparsity constraints also encourage this condition to be satisfied. Based on Bello et al. (2022), we implemented a path-following shown in Algorithm 2.

The optimization model (17) is observed well for linear Gaussian SEMs with equal variance as well as other equal variance SEMs. Meanwhile, for unequal variance, or normalized data from linear Gaussian SEMs with equal variance where the scale information are missing, MSE score function is not consistent and often provides misleading information about the underlying DAG. Additionally, as observed by Ng et al. (2023), the initialization of adjacency matrices in cases of unequal variance can significantly affect performance, suggesting that non-convexity may pose a serious challenge in such scenarios.

4 NON-CONVEXITY ANALYSIS OF ANALYTIC DAG CONSTRAINTS

The non-convexity of a function can be analyzed through the analysis of its Hessian. Particularly for our analytic DAG constraints, its Hessian can be obtained using the following proposition and then the non-convexity can be analyzed by analysis the spectral radius of Hessian.

Proposition 6. *The Hessian of DAG constraints (6) can be obtained as follows:*

$$\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}}) = \mathbf{K}_{dd} \sum_{i=2}^{\infty} i c_i \sum_{j=0}^{i-2} [\tilde{\mathbf{B}}^j]^{\top} \otimes [\tilde{\mathbf{B}}^{i-2-j}], \quad (18)$$

where \otimes denotes the Kronecker product, and $\mathbf{K}_{dd} \in \{0, 1\}^{d^2 \times d^2}$ is the commutation matrix satisfies that for any $d \times d$ matrix \mathbf{A}

$$\mathbf{K}_{d,d} \text{vec}(\mathbf{A}) = \text{vec}(\mathbf{A}^{\top}).$$

¹the constant $c_0 d$ can be dropped because $\text{tr} f(\mathbf{B} \odot \mathbf{B})$ is bounded below by $c_0 d$, detailed derivation is provided in the supplementary file.

Obviously, the Hessian Equation (18) is symmetric and not positive semi-definite. One widely used way to convexify Hessian is to find a positive scalar η such that

$$\Delta = \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}}) + \eta \mathbb{I}, \quad (19)$$

becomes positive semi-definite. It require η to be no less than the absolute value of the most negative eigenvalue of $\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}})$. Here the Hessian are symmetric matrix with all non-negative entries. For this kind of matrices the absolute value of the most negative eigenvalue of $\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}})$ is upper bounded by the spectral radius of Hessian, and the bound is tight under certain conditions (Spielman, 2012). Thus it would be nature to use the spectral radius of Hessian to measure the level of non-convexity of the analytic DAG constraints.

The Hessian $\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}})$ can be viewed as linear combinations of a series of symmetric matrices $i \mathbf{K}_{dd} \sum_{j=0}^{i-2} [\tilde{\mathbf{B}}^j]^\top \otimes [\tilde{\mathbf{B}}^{i-2-j}]$ with all non-negative entries. The commutation matrix \mathbf{K}_{dd} (Magnus and Neudecker, 1979) would not have any effects on the spectral radius as it is orthonormal. Thus larger c_i would result the spectral radius of a single term $i c_i \mathbf{K}_{dd} \sum_{j=0}^{i-2} [\tilde{\mathbf{B}}^j]^\top \otimes [\tilde{\mathbf{B}}^{i-2-j}]$ to increase, and finally lead the spectral radius of the Hessian to increase as the following proposition.

Proposition 7. For two analytic function $f_1(x) = c_{0,1} + \sum_{i=1}^{\infty} c_{i,1} x^i$ and $f_2(x) = c_{0,2} + \sum_{i=1}^{\infty} c_{i,2} x^i$, if $\forall i \geq 1$ we have $c_{i,1} \geq c_{i,2} > 0$, then

$$\rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_1(\tilde{\mathbf{B}})) \geq \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_2(\tilde{\mathbf{B}})), \quad (20)$$

where $\rho(\cdot)$ denotes the spectral radius of a matrix.

Proposition 7 suggests that the spectral radius of the Hessian would increase if the coefficients c_i in the analytic function increases. This implies that DAG constraints with larger c_i may gain benefits from gradient vanishing, but suffers from non-convexity. In fact, using Proposition 7 it would be straightforward to get the following corollary, which provides the level of non-convexity comparison for several DAG constraints.

Corollary 8.

$$\rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} \exp(\tilde{\mathbf{B}})) \leq \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_{\log}^s(\tilde{\mathbf{B}})) \leq \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_{inv}^s(\tilde{\mathbf{B}})). \quad (21)$$

The optimization problem (17) can be viewed as a convex objective plus one non-convex constraint, and the convex mean square error (MSE) loss may play different roles in different scenarios. For data with known scale, the MSE loss is consistent and thus it provides enough information to identify the underlying model and thus the non-convexity may not be a serious issue. This is because in the path-following optimization framework (provided in Algorithm 2), at the beginning the optimization direction are dominated by the MSE loss so that it will push the candidates to a point that is not far from global optimal. Thus DAG constraints with finite convergence radius is preferred to escape from gradient vanishing. Meanwhile, for DAG learning problem with unknown scale, the MSE loss may not be very informative to the underlying graph structure. In this case, the highly non-convex DAG constraints may lead to the optimizer to get trapped into a local minimum easily, and thus we may need additional constraints to reduce the search space, which may possibly make the objective flatter. In our experiments, we find that by allowing only edges to exist between nodes with strong correlation can significantly improve the performance.

5 EXPERIMENTS

In the experiment, we compared the performance of different analytic DAG constraints in the same path-following optimization framework. We implemented the path-following algorithm (provided in Algorithm 2) using PyTorch (Paszke et al., 2019) based on the path-following optimizer in Bello et al. (2022). For analytic DAG constraints with infinite convergence radius, we consider the exponential-based DAG constraints. For analytic DAG constraints with finite convergence radius, we consider the following 4 different DAG constraints generated by the differentiation operator or multiply operator:

- Order-1: $\text{tr} f_{\log}^s(\mathbf{B} \odot \mathbf{B}) = 0$;

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Graphs	#Nodes	DAGMA	Order 1	Order 2	Order 3	Order 4
ER2-Gaussian	500	44.90 ± 32.95	33.40 ± 23.46	31.70 ± 19.47	30.60 ± 19.07	29.80 ± 20.97
	1000	94.80 ± 35.80	69.60 ± 27.64	55.60 ± 19.13	52.40 ± 19.86	57.30 ± 21.39
	2000	235.40 ± 62.76	176.00 ± 47.77	153.30 ± 35.92	135.60 ± 38.91	131.00 ± 28.65
ER3-Gaussian	500	125.30 ± 44.55	101.10 ± 39.03	90.30 ± 39.56	93.90 ± 31.26	92.90 ± 45.56
	1000	339.60 ± 67.80	242.80 ± 72.21	210.30 ± 60.98	184.90 ± 47.44	165.80 ± 35.95
	2000	669.50 ± 140.61	610.70 ± 136.84	555.30 ± 106.01	479.50 ± 88.72	424.90 ± 64.39
ER4-Gaussian	500	307.60 ± 116.53	261.40 ± 102.81	263.00 ± 122.34	246.70 ± 110.28	223.80 ± 97.46
	1000	878.50 ± 174.96	689.20 ± 165.62	695.50 ± 134.41	619.80 ± 150.43	626.30 ± 157.59
	2000	1922.30 ± 187.69	1785.40 ± 184.47	1779.30 ± 211.23	1655.40 ± 181.75	1574.10 ± 152.23
ER2-Exp	500	58.20 ± 31.58	40.50 ± 26.93	28.90 ± 16.60	31.00 ± 25.67	35.20 ± 34.32
	1000	93.90 ± 33.96	68.70 ± 23.20	54.00 ± 16.26	50.90 ± 17.29	57.90 ± 24.93
ER3-Exp	500	142.70 ± 50.13	106.00 ± 39.15	95.10 ± 32.68	99.60 ± 37.94	100.30 ± 47.52
	1000	321.10 ± 83.82	242.80 ± 68.51	212.40 ± 67.87	187.60 ± 61.87	173.10 ± 49.03
ER4-Exp	500	336.00 ± 124.19	292.70 ± 123.41	294.90 ± 130.66	254.40 ± 133.05	214.70 ± 84.29
	1000	879.40 ± 162.98	718.20 ± 127.12	710.60 ± 151.41	640.50 ± 148.24	619.70 ± 133.44
ER2-Gumbel	500	45.10 ± 33.28	22.60 ± 20.04	21.30 ± 18.41	19.80 ± 16.39	16.20 ± 11.91
	1000	80.50 ± 42.65	49.90 ± 24.54	39.90 ± 14.04	36.80 ± 15.18	45.90 ± 23.18
ER3-Gumbel	500	147.10 ± 54.19	94.10 ± 40.87	76.60 ± 60.77	60.60 ± 31.34	85.30 ± 50.75
	1000	297.90 ± 72.40	215.40 ± 52.35	185.00 ± 71.98	173.90 ± 57.09	147.70 ± 40.51
ER4-Gumbel	500	338.80 ± 127.56	273.70 ± 131.13	257.50 ± 111.06	232.40 ± 121.98	234.70 ± 149.59
	1000	919.90 ± 182.38	722.80 ± 177.86	734.80 ± 177.78	620.70 ± 187.56	564.10 ± 170.46

Table 1: DAG learning performance (measured in structural hamming distance, the lower the better, best results in **bold**) of different algorithms on ER{2,3,4} graphs with different noise distributions. All our algorithms performs better than the previous state-of-the-arts DAGMA (Bello et al., 2022), and as higher order DAG constraints suffers less to gradient vanishing, it tends to have better performance.

Graphs	#Nodes	DAGMA(Bello et al., 2022)	Order 1	Order 2	Order 3	Order 4
SF2	500	31.40 ± 43.51	24.30 ± 43.90	32.40 ± 49.38	34.20 ± 45.56	41.50 ± 48.45
	1000	44.90 ± 34.38	41.20 ± 36.02	22.50 ± 13.21	29.20 ± 20.07	58.10 ± 27.58
	2000	189.80 ± 99.47	162.90 ± 73.30	172.10 ± 74.35	152.20 ± 90.29	172.60 ± 124.12
SF3	500	58.10 ± 33.90	51.10 ± 32.10	41.10 ± 17.91	49.80 ± 24.58	71.00 ± 23.46
	1000	169.40 ± 60.82	158.10 ± 46.70	161.20 ± 55.25	162.50 ± 57.54	195.40 ± 75.60
	2000	928.70 ± 148.70	896.00 ± 101.85	897.10 ± 146.78	891.50 ± 143.40	999.70 ± 206.36
SF4	500	131.20 ± 42.63	136.80 ± 41.71	134.40 ± 39.34	128.90 ± 36.68	151.60 ± 37.05
	1000	431.70 ± 119.22	404.00 ± 88.89	400.30 ± 76.49	386.90 ± 93.16	394.50 ± 111.57
	2000	1525.10 ± 299.02	1500.50 ± 297.88	1444.70 ± 291.45	1395.60 ± 264.90	1418.90 ± 228.86
SF2	500	25.90 ± 44.45	23.40 ± 44.41	32.10 ± 49.08	35.00 ± 48.84	37.20 ± 45.21
	1000	43.70 ± 34.48	41.20 ± 36.02	32.00 ± 34.13	29.10 ± 19.68	59.10 ± 30.34
SF3	500	57.70 ± 33.68	57.70 ± 33.64	41.80 ± 20.37	43.20 ± 15.75	66.70 ± 24.36
	1000	177.10 ± 67.53	165.40 ± 57.70	171.60 ± 66.80	175.10 ± 69.09	195.90 ± 80.37
SF4	500	127.50 ± 40.84	132.80 ± 40.39	129.90 ± 42.07	135.50 ± 44.21	152.40 ± 39.94
	1000	408.80 ± 119.71	419.70 ± 108.01	388.50 ± 53.01	394.30 ± 88.95	395.30 ± 109.37
SF2	500	23.10 ± 44.78	17.70 ± 44.51	16.70 ± 44.15	20.00 ± 45.75	33.40 ± 49.30
	1000	29.20 ± 24.77	24.70 ± 24.85	12.50 ± 11.40	16.20 ± 13.96	47.90 ± 25.69
SF3	500	33.50 ± 32.98	25.20 ± 27.85	19.40 ± 12.37	19.00 ± 7.44	50.00 ± 22.41
	1000	107.50 ± 50.50	114.50 ± 59.80	106.60 ± 64.77	103.70 ± 58.15	133.60 ± 88.43
SF4	500	77.70 ± 41.43	76.20 ± 41.86	67.90 ± 26.47	79.20 ± 23.87	101.40 ± 22.37
	1000	333.10 ± 118.06	348.80 ± 110.93	309.20 ± 51.86	321.50 ± 83.13	339.70 ± 111.17

Table 2: DAG learning performance (measured in structural hamming distance, the lower the better, best results in **bold**) of different algorithms on SF{2,3,4} graphs with different noise distributions. Our algorithms usually performs better than the previous state-of-the-arts DAGMA(Bello et al., 2022).

- Order-2: $\text{tr}[f_{inv}^s(\mathbf{B} \odot \mathbf{B} / s)] = d$;
- Order-3: $\text{tr}[f_{inv}^s(\mathbf{B} \odot \mathbf{B} / s)]^2 = d$;
- Order-4: $\text{tr}[f_{inv}^s(\mathbf{B} \odot \mathbf{B} / s)]^3 = d$.

In our experiments, we use the same annealing strategy for s as Bello et al. (2022). During the optimization, the spectral radius of $\mathbf{B} \odot \mathbf{B}$ may be larger than s , which make the DAG constraints

	PC	GES	DAGMA	Exponential	Order 1	Order 2	Order 3	Order 4
SHD	563.9 ± 23.84	4490.2 ± 62.52	588.8 ± 18.33	488.6 ± 24.29	429.6 ± 24.73	410.6 ± 15.25	401.0 ± 16.64	389.4 ± 16.70

Table 3: DAG learning performance (measured in structural hamming distance, the lower the better, best results in **bold**) of different algorithms on 1000-node ER1 graphs with Gaussian noise with observation data normalized. Our algorithm performs better than the previous approaches, and as higher order DAG constraints suffers less to gradient vanishing, it tends to have better performance. We compare differential DAG learning approaches with conditional independent test based PC (Spirtes and Glymour, 1991) algorithm and score based GES (Chickering, 2002) algorithm. The result is reported in the format of average \pm standard derivation gathered from 10 different simulations.

invalid. In this case, we use the strategy provided in Algorithm 1 to reset s . We also tried the DAG constraints (Zhang et al., 2022) with their code provided in their appendix, but it do have some numeric issues for large scale graphs.

We compare the performance of these DAG constraints using two different settings: linear SEM with known ground truth scale and with unknown ground truth scale. We also compare these methods with constraint based PC (Spirtes and Glymour, 1991) algorithm and score based combinatorial search algorithm GES (Chickering, 2002) implemented by Kalainathan et al. (2020).

5.1 LINEAR SEM WITH KNOWN GROUND TRUTH SCALE

For linear SEM with a known ground truth scale, our experimental setting is similar to Bello et al. (2022); Zhang et al. (2022); Zheng et al. (2018). We generated two different types of random graphs: ER (Erdős-Rényi) and SF (Scale-Free) graphs with different numbers of expected edges. We use ER_n (SF_n) to denote graphs with d nodes and nd expected edges. Edge weights generated from a uniform distribution over the union of two intervals $[-2, -0.5] \cup [0.5, 2.0]$ are assigned to each edge to form a weighted adjacency matrix \mathbf{B} . Then, n samples are generated from the linear SEM $\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{e}$ to form an $n \times d$ data matrix \mathbf{X} , where the noise \mathbf{e} is iid sampled from Gaussian, Exponential, or Gumbel distribution. As Bello et al. (2022); Zhang et al. (2022); Zheng et al. (2018) achieved nearly perfect results on small and sparse graphs, we considered more challenging large and denser graphs in our experiments. We set the sample size $n = 1000$ and consider 3 different numbers of nodes $d = 500, 1000, 2000$. For each setting, we conducted 10 random simulations to obtain an average performance. All these experiments were performed using an A100 GPU, and all computations were done in double precision. Our algorithms were compared with the previous state-of-the-art approach DAGMA Bello et al. (2022). The original version of DAGMA Bello et al. (2022) used numpy and ran on CPU; we replaced numpy with cupy to get a GPU version of DAGMA, which performed identically to the CPU version.

The results on ER2, ER3, and ER4 graphs are shown in Table 1. In all cases, our algorithms outperformed the previous state-of-the-art DAGMA. Our Order-1 algorithm is very similar to DAGMA, except for our annealing strategy of s derived from our theory, which indicates the efficiency of our theory. Furthermore, our theory shows that higher-order constraints suffer less from gradient vanishing, and in the experimental results, we observed that the performance of higher-order DAG constraints outperformed lower-order ones in most cases. The results of SF2, SF3, and SF4 graphs are shown in Table 2. On scale-free graphs, our algorithms usually performed better than DAGMA, and the higher-order constraints, Order-2 and Order-3, often outperformed Order-1. The performance of Order-4 constraints was not good, possibly due to stronger non-convexity.

The DAGMA algorithm actually employed the same DAG constraints as our Order-1 method, but with a different strategy to search for s . Our search strategy, derived from properties of analytic functions, provides a tight bound for s , allowing a smaller s to be used than DAGMA without sacrificing the numeric stability. As a result, our algorithm suffers less from gradient vanishing and achieves better performance.

In terms of running time, all algorithms had similar running times, typically about 5 minutes for a 500-node graph, 10-20 minutes for a 1000-node graph, and around 2 hours for a 2000-node graph. Due to limited time and resources, we only considered $d = 2000$ for Gaussian noises, and for other cases, we only considered $d = 500, 1000$.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

	Original GRAN	Order 1	Order 2	Order 3	Order 4
SHD	15	13	13	13	13
SHD-CPDAG	10	9	9	9	9
SHD	13	13	13	13	13
SHD-CPDAG	10	9	9	9	9

Table 4: Nonlinear DAG learning performance (measured in structural hamming distance on DAG ad CPDAG, the lower the better) of different DAG constraints on Sachs et al. (2005)’s dataset . Different DAG constraints was plugged into the GRAN (Lachapelle et al., 2020) framework. **Top two rows**: results obtained from single-precision mode. **Bottom two rows**: results obtained from double-precision mode.

5.2 LINEAR SEM WITH UNKNOWN GROUND TRUTH SCALE

For linear SEM with an unknown ground truth scale, we applied the same data generation process as for the linear SEM with a known ground truth scale and Gaussian noise, but normalized the generated data \mathbf{X} to have zero mean and unit variance. In this normalization procedure, the scale information of the variables is removed from the data. Particularly for Gaussian noise, in this case, the true DAG is not identifiable, and we may only identify it to a Markov Equivalent Class. Previously, it has been observed that direct optimization over (17) may result in poor performance, mainly due to the non-convex nature. In our experiments, we added an additional constraint to only allow edges to exist between highly correlated nodes. We first computed the Pearson correlation coefficients between every pair of nodes, and if the absolute value of the coefficient between two nodes is larger than 0.1, then we allow the edge to exist. During optimization, at every gradient descent step, we removed the disallowed edges from the candidate graph. We generated 10 instances of 1000-node ER1 graphs with Gaussian noise, and 1000 observational samples were generated for each instance.

The results are shown in Table 3. Our algorithm outperforms PC (Spirtes and Glymour, 1991) and GES (Chickering, 2002) in terms of SHD. Although higher-order DAG constraints may suffer more from non-convexity, by adding proper constraints on the candidate graphs, we can still achieve satisfactory results. In the results, we can see that the performance of higher-order constraints is better than that of lower-order ones, and also better than the exponential-based DAG constraints, which suggests that gradient vanishing may still be one important reason for poor performance.

5.3 EXPERIMENTAL RESULTS ON NONLINEAR CASES

Our DAG constraints can also be extended to continuous nonlinear DAG learning approaches by replacing their original DAG constraints. We incorporated our DAG constraints into Lachapelle et al. (2020) to model nonlinear Structural Equation Models (SEMs) and conducted experiments using Sachs et al. (2005)’s dataset pre-processed by Lachapelle et al. (2020)². The GraN-DAG algorithm can operate in both single-precision and double-precision modes. The experimental results are shown in Table 4. The results suggest that DAG constraints with a finite spectral radius suffer less from gradient vanishing and, consequently, from numerical truncation errors. In contrast, the original GraN-DAG algorithm experiences gradient vanishing, particularly when running in single-precision mode, as higher-order constraints that prevent long loops are truncated due to limited machine precision.

6 CONCLUSION

The continuous differentiable DAG constraints play an important role in the continuous DAG learning algorithms. We show that many of these DAG constraints can be formulated using analytic functions. Several functional operators, including differentiation, summation, and multiplication, can be leveraged to create novel DAG constraints based on existing ones. Using these properties, we designed a series of DAG constraints and designed an efficient algorithm to evaluate these DAG constraints. Experiments on various settings show that our DAG constraints outperform previous state-of-the-arts approaches.

²Available at <https://github.com/kurowasan/GraN-DAG>.

540 REFERENCES

541 Kevin Bello, Bryon Aragam, and Pradeep Kumar Ravikumar. Dagma: Learning dags via m-matrices and a

542 log-determinant acyclicity characterization. In *Advances in Neural Information Processing Systems*, 2022.

543

544 David Maxwell Chickering. Learning Bayesian networks is NP-complete. In *Learning from Data: Artificial*

545 *Intelligence and Statistics V*. Springer, 1996.

546 David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning*

547 *research*, 3(Nov):507–554, 2002.

548

549 Max Chickering, David Heckerman, and Chris Meek. Large-sample learning of bayesian networks is np-hard.

550 *Journal of Machine Learning Research*, 5, 2004.

551 Diego Colombo, Marloes H Maathuis, Markus Kalisch, and Thomas S Richardson. Learning high-dimensional

552 directed acyclic graphs with latent and selection variables. *The Annals of Statistics*, pages 294–321, 2012.

553 James Cussens. Bayesian network learning with cutting planes. In *Conference on Uncertainty in Artificial*

554 *Intelligence*, 2011.

555

556 Nicholas J Higham. *Functions of matrices: theory and computation*. SIAM, 2008.

557 Diviyani Kalainathan, Olivier Goudet, and Ritik Dutta. Causal discovery toolbox: Uncovering causal relationships

558 in python. *The Journal of Machine Learning Research*, 21(1):1406–1410, 2020.

559 Mikko Koivisto and Kismat Sood. Exact Bayesian structure discovery in Bayesian networks. *Journal of Machine*

560 *Learning Research*, 5(Dec):549–573, 2004.

561 Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.

562 Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural

563 DAG learning. In *International Conference on Learning Representations*, 2020.

564

565 Po Ling Loh and Peter Bühlmann. High-dimensional learning of linear causal networks via inverse covariance

566 estimation. *Journal of Machine Learning Research*, 2013.

567

568 Jan R Magnus and Heinz Neudecker. The commutation matrix: some properties and applications. *The annals of*

569 *statistics*, 7(2):381–394, 1979.

570 Jan R Magnus and Heinz Neudecker. *Matrix differential calculus with applications in statistics and econometrics*.

571 John Wiley & Sons, 2019.

572 Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and DAG constraints for learning

573 linear DAGs. *Advances in Neural Information Processing Systems*, 33, 2020.

574

575 Ignavier Ng, Sébastien Lachapelle, Nan Rosemary Ke, Simon Lacoste-Julien, and Kun Zhang. On the conver-

576 gence of continuous constrained optimization for structure learning. In *International Conference on Artificial*

577 *Intelligence and Statistics*, pages 8176–8198. PMLR, 2022.

578 Ignavier Ng, Biwei Huang, and Kun Zhang. Structure learning with continuous optimization: A sober look and

579 beyond. *arXiv preprint arXiv:2304.02146*, 2023.

580 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,

581 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep

582 learning library. *Advances in neural information processing systems*, 32, 2019.

583

584 Judea Pearl. *Models, reasoning and inference*. Cambridge, UK: Cambridge University Press, 19, 2000.

585 Jonas Peters and Peter Bühlmann. Identifiability of Gaussian structural equation models with equal error

586 variances. *Biometrika*, 101(1):219–228, 2013.

587 Garvesh Raskutti and Caroline Uhler. Learning directed acyclic graph models based on sparsest permutations.

588 *Stat*, 7(1):e183, 2018.

589

590 Alexander Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag! causal discovery

591 benchmarks may be easy to game. *Advances in Neural Information Processing Systems*, 34:27772–27784,

592 2021.

593 Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal protein-signaling

networks derived from multiparameter single-cell data. *Science*, 308(5721):523–529, 2005.

594 Rajen D Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance
595 measure. *The Annals of Statistics*, 48(3):1514–1538, 2020.

596 Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-Gaussian acyclic model
597 for causal discovery. *Journal of Machine Learning Research*, 7(Oct):2003–2030, 2006.

598 Tomi Silander and Petri Myllymäki. A simple approach for finding the globally optimal Bayesian network
599 structure. In *Conference on Uncertainty in Artificial Intelligence*, 2006.

600 Ajit P. Singh and Andrew W. Moore. Finding optimal Bayesian networks by dynamic programming. Technical
601 report, Carnegie Mellon University, 2005.

602 Daniel A. Spielman. Lecture notes in spectral graph theory: The adjacency matrix and the n th eigenvalue, Septem-
603 ber 2012. URL <https://www.cs.yale.edu/homes/spielman/561/2012/lect03-12.pdf>.

604 Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social science
605 computer review*, 9(1):62–72, 1991.

606 Peter Spirtes, Chris Meek, and Thomas Richardson. Causal inference in the presence of latent variables and
607 selection bias. In *Conference on Uncertainty in Artificial Intelligence*, 1995.

608 Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*.
609 MIT press, 2000.

610 Matthew J Vowels, Necati Cihan Camgoz, and Richard Bowden. D’ya like dags? a survey on structure learning
611 and causal discovery. *arXiv preprint arXiv:2103.02582*, 2021.

612 Dennis Wei, Tian Gao, and Yue Yu. Dags with No Fears: A closer look at continuous optimization for learning
613 bayesian networks. *Advances in Neural Information Processing Systems*, 33:3895–3906, 2020.

614 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. DAG-GNN: Dag structure learning with graph neural networks. In
615 *International Conference on Machine Learning*, pages 7154–7163. PMLR, 2019.

616 Yue Yu, Tian Gao, Naiyu Yin, and Qiang Ji. DAGs with no curl: An efficient dag structure learning approach. In
617 *Proceedings of the 38th International Conference on Machine Learning*, 2021.

618 Changhe Yuan and Brandon Malone. Learning optimal Bayesian networks: A shortest path perspective. *Journal
619 of Artificial Intelligence Research*, 48(1):23–65, 2013.

620 Changhe Yuan, Brandon Malone, and Xiaojian Wu. Learning optimal Bayesian networks using A* search. In
621 *International Joint Conference on Artificial Intelligence*, 2011.

622 Bin Zhang, Chris Gaiteri, Liviu-Gabriel Bodea, Zhi Wang, Joshua McElwee, Alexei A Podtelezchnikov, Chun-
623 sheng Zhang, Tao Xie, Linh Tran, Radu Dobrin, et al. Integrated systems approach identifies genetic nodes
624 and networks in late-onset alzheimer’s disease. *Cell*, 153(3):707–720, 2013.

625 Zhen Zhang, Ignavier Ng, Dong Gong, Yuhang Liu, Ehsan Abbasnejad, Mingming Gong, Kun Zhang, and
626 Javen Qinfeng Shi. Truncated matrix power iteration for differentiable dag learning. *Advances in Neural
627 Information Processing Systems*, 35:18390–18402, 2022.

628 Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. DAGs with No Tears: Continuous
629 optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.

630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Appendices

A PROOF OF PROPOSITIONS

Our proof are also based on the well-known properties of analytic functions listed as follows:

1. Let $f_1(x), f_2(x)$ be analytic functions on $(-r_1, r_1)$ and $(-r_2, r_2)$, then $f_1(x) + f_2(x)$ and $f_1(x)f_2(x)$ are analytic functions on $(-\min(r_1, r_2), \min(r_1, r_2))$;
2. Let $f(x)$ be analytic function on $(-r, r)$, then $\partial f(x)/\partial x$ is an analytic function on $(-r, r)$.

A.1 LEMMAS REQUIRED FOR PROOFS

Lemma 9. Let $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ be the weighted adjacency matrix of a graph \mathcal{G} with d vertices, \mathcal{G} is a DAG if and only if $\tilde{\mathbf{B}}^d = \mathbf{0}$.

Proof. See Proposition 3.1 of Zhang et al. (2022). □

Lemma 10. Let $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ be the weighted adjacency matrix of a graph \mathcal{G} with d vertices, \mathcal{G} is a DAG if and only

$$\text{tr}\left(\sum_{i=1}^d c_i \tilde{\mathbf{B}}^i\right) = 0,$$

where $c_i > 0 \forall i$.

Proof. See Wei et al. (2020). □

A.2 PROOF OF PROPOSITION 1

Proposition 1. Let $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with $\rho(\tilde{\mathbf{B}}) \leq r$ be the weighted adjacency matrix of a directed graph \mathcal{G} , and let f be an analytic function in the form of (3), where we further assume $\forall i > 0$ we have $c_i > 0$, then \mathcal{G} is acyclic if and only if

$$\text{tr}\left[f(\tilde{\mathbf{B}})\right] = c_0 d.$$

Proof. Without loss of generality, assume that f can be formulated as:

$$f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i; \forall i, c_i > 0; \lim_{i \rightarrow \infty} c_i / c_{i+1} > 0. \quad (22)$$

First if \mathcal{G} is acyclic, by Lemma 9 we must have

$$\tilde{\mathbf{B}}^k = \mathbf{0} \forall k \geq d, \quad (23)$$

which also indicates that $\rho(\tilde{\mathbf{B}}) = 0$. Thus we have

$$\begin{aligned} \text{tr}\left[f(\tilde{\mathbf{B}})\right] &= \text{tr}\left[c_0 \mathbf{I} + \sum_{i=1}^d c_i \tilde{\mathbf{B}}^i + \underbrace{\sum_{i=d+1}^{\infty} c_i \tilde{\mathbf{B}}^i}_{\text{Equals 0, By Lemma 9}}\right] \\ &= \text{tr}[c_0 \mathbf{I}] + \underbrace{\text{tr}\left[\sum_{i=1}^d c_i \tilde{\mathbf{B}}^i\right]}_{\text{Equals 0, By Lemma 10}} \\ &= c_0 d. \end{aligned} \quad (24)$$

On the other hand, if $\text{tr} [f(\tilde{\mathbf{B}})] = c_0 d$, we must have that

$$\text{tr} \left[\sum_{i=1}^{\infty} c_i \tilde{\mathbf{B}}^i \right] = 0.$$

By the fact all entries of $\tilde{\mathbf{B}}$ are positive, we have that

$$0 \leq \text{tr} \left[\sum_{i=1}^d c_i \tilde{\mathbf{B}}^i \right] \leq \left[\sum_{i=1}^{\infty} c_i \tilde{\mathbf{B}}^i \right] = 0. \quad (25)$$

Then we must have

$$\text{tr} \left[\sum_{i=1}^d c_i \tilde{\mathbf{B}}^i \right] = 0.$$

Finally by Lemma 10 we have that \mathcal{G} is a DAG. \square

A.3 PROOF OF PROPOSITION 2

In all the paper, we consider analytic functions f from the functional class \mathcal{F} defined in (5).

Proposition 2. *There exists some real number r , where for all $\{\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} \mid \rho(\tilde{\mathbf{B}}) < r\}$, the derivative of $\text{tr} [f(\tilde{\mathbf{B}})]$ w.r.t. $\tilde{\mathbf{B}}$ is*

$$\nabla_{\tilde{\mathbf{B}}} \text{tr} [f(\tilde{\mathbf{B}})] = [\nabla_x f(x)|_{x=\tilde{\mathbf{B}}}]^{\top}.$$

Proof. Without loss of generality, assume that f can be formulated as:

$$f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i; \forall i, c_i > 0; \lim_{i \rightarrow \infty} c_i / c_{i+1} > 0. \quad (26)$$

For some i by basic matrix differentiation we have

$$\frac{\partial \text{tr} \tilde{\mathbf{B}}^i}{\partial \tilde{\mathbf{B}}} = (i \tilde{\mathbf{B}}^{i-1})^{\top}, \quad (27)$$

and then by the properties of power series we have

$$\begin{aligned} \nabla_{\tilde{\mathbf{B}}} \text{tr} [f(\tilde{\mathbf{B}})] &= \nabla_{\tilde{\mathbf{B}}} \text{tr} \left[c_0 \mathbf{I} + \sum_{i=1}^{\infty} c_i \tilde{\mathbf{B}}^i \right] \\ &= \sum_{i=1}^{\infty} \nabla_{\tilde{\mathbf{B}}} \text{tr} c_i \tilde{\mathbf{B}}^i \\ &= \left[\sum_{i=1}^{\infty} c_i i \tilde{\mathbf{B}}^{i-1} \right]^{\top} \\ &= \left[\sum_{i=1}^{\infty} c_i i x^{i-1} \Big|_{x=\tilde{\mathbf{B}}} \right]^{\top} = [\nabla_x f(x)|_{x=\tilde{\mathbf{B}}}]^{\top}, \end{aligned} \quad (28)$$

where we can exchange ∇ and $\sum_{i=1}^{\infty}$ because after the exchanging the new power series will still converge (by properties of analytic functions). \square

A.4 PROOF OF PROPOSITION 3

Proposition 3. *Let $f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i \in \mathcal{F}$ be a analytic function on $(-r, r)$, and let n be arbitrary integer larger than 1, then $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with spectral radius $\rho(\tilde{\mathbf{B}}) \leq r$ forms a DAG if and only if*

$$\text{tr} \left[\frac{\partial^n f(x)}{\partial x^n} \Big|_{x=\tilde{\mathbf{B}}} \right] = n! c_n.$$

756 *Proof.* By properties of analytic functions, the n^{th} order derivative of an analytic function $f(x)$ on
757 $(-r, r)$ is still an analytic function on $(-r, r)$. Particularly for $f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i \in \mathcal{F}$, we have
758

$$\begin{aligned}
759 \quad \frac{\partial^n f(x)}{\partial x^n} &= \sum_{i=1}^{\infty} \frac{\partial^n c_i x^i}{\partial x^n} \\
760 &= \sum_{i=n}^{\infty} \frac{\partial^n c_i x^i}{\partial x^n} \\
761 &= \sum_{i=n}^{\infty} \left[c_i x^{i-n} \prod_{k=i-n+1}^n k \right] \\
762 &= n! c_n + \sum_{i=1}^{\infty} \left[c_{i+n} x^i \prod_{k=i}^{n+i} k \right], \tag{29} \\
763 & \\
764 & \\
765 & \\
766 & \\
767 & \\
768 & \\
769 & \\
770 &
\end{aligned}$$

771 where by the fact $c_i > 0 \forall i > 1$, we have that $\frac{\partial^n f(x)}{\partial x^n} \in \mathcal{F}$. Then by Proposition 1 we immediately
772 proved the proposition. \square
773

774 A.5 PROOF OF PROPOSITION 4

775 **Proposition 4.** Let $f_1(x) = c_0^1 + \sum_{i=1}^{\infty} c_i^1 x^i \in \mathcal{F}$, and $f_2(x) = c_0^2 + \sum_{i=1}^{\infty} c_i^2 x^i \in \mathcal{F}$. Then for an ad-
776 jacency matrix $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with spectral radius $\rho(\tilde{\mathbf{B}}) \leq \min(\lim_{i \rightarrow \infty} c_i^1 / c_{i+1}^1, \lim_{i \rightarrow \infty} c_i^2 / c_{i+1}^2)$,
777 the following three statements are equivalent:
778

- 779 1. $\tilde{\mathbf{B}}$ forms a DAG;
- 780 2. $\text{tr}[f_1(\tilde{\mathbf{B}}) + f_2(\tilde{\mathbf{B}})] = (c_0^1 + c_0^2)d$;
- 781 3. $\text{tr}[f_1(\tilde{\mathbf{B}})f_2(\tilde{\mathbf{B}})] = c_0^1 c_0^2 d$.

782 *Proof.* By properties of analytic functions, we have
783

$$784 \quad f_1(x) + f_2(x) = c_0^1 + c_0^2 + \sum_{i=1}^{\infty} (c_i^1 + c_i^2) x^i \tag{30}$$

785 is an analytic function, and its convergence radius is given by
786

$$787 \quad \lim_{i \rightarrow \infty} (c_i^1 + c_i^2) / (c_{i+1}^1 + c_{i+1}^2) = \min(\lim_{i \rightarrow \infty} c_i^1 / c_{i+1}^1, \lim_{i \rightarrow \infty} c_i^2 / c_{i+1}^2), \tag{31}$$

788 and thus by Proposition 1 the statement 1 and 2 are equivalent. Similarly by properties of analytic
789 functions statement 1 and 3 are equivalent. Thus the 3 statements are equivalent. \square
790

791 A.6 PROOF OF PROPOSITION 5

792 **Proposition 5.** Let n be any positive integer, the adjacency matrix $\tilde{\mathbf{B}} \in \{\hat{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\hat{\mathbf{B}}) \leq s\}$ if
793 and only if
794

$$795 \quad \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n} = d,$$

796 and the gradients of the DAG constraints satisfies that $\forall \tilde{\mathbf{B}} \in \{\hat{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\hat{\mathbf{B}}) \leq s\}$
797

$$798 \quad \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n}\| \leq \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n-k}\|,$$

799 where k is an arbitrary positive integer, and $\|\cdot\|$ denote an arbitrary matrix norm induced by vector
800 p -norm.
801

802 *Proof.* By Proposition 4 or Proposition 3, it would be straightforward that $\text{tr}(\mathbf{I} - \tilde{\mathbf{B}})^{-n} = d$ is a
803 necessary and sufficient condition for an adjacency matrix $\tilde{\mathbf{B}} \in \{\hat{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\hat{\mathbf{B}}) \leq s\}$ to form a
804 DAG.
805
806
807
808
809

810 For the norm of gradients, it is straightforward that

$$811 \frac{\partial(1-x)^{-n}}{\partial x} = n(1-x)^{-n-1}. \quad (32)$$

812 For arbitrary n we have

$$813 (1-x)^{-n} = 1 + \sum_{i=1}^{\infty} \left[\prod_{j=n}^{n+i-1} j \right] x^i, \quad (33)$$

814 and obviously the coefficients is monotonic increasing w.r.t. n . Thus by the fact $\forall \tilde{\mathbf{B}} \in \{\hat{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\hat{\mathbf{B}}) \leq s\}$ we have for any $j > 0, k > 0$

$$815 \|\mathbf{I} - \tilde{\mathbf{B}}^{-j}\| \leq \|\mathbf{I} - \tilde{\mathbf{B}}^{-j-k}\|. \quad (34)$$

816 As a result, we have

$$817 \begin{aligned} \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}})^{-n}\| &= n \|\mathbf{I} - \tilde{\mathbf{B}}\|^{-n-1} \leq (n+k) \|\mathbf{I} - \tilde{\mathbf{B}}\|^{-n-1} \\ &\leq (n+k) \|\mathbf{I} - \tilde{\mathbf{B}}\|^{-n-k-1} = \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}})^{-n-k}\|. \end{aligned} \quad (35)$$

818 □

819 A.7 PROOF OF PROPOSITION 6

820 **Proposition 6.** *Thus the Hessian of DAG constraints (6) can be obtained as follows:*

$$821 \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}}) = \mathbf{K}_{dd} \sum_{i=2}^{\infty} i c_i \sum_{j=0}^{i-2} [\tilde{\mathbf{B}}^j]^\top \otimes [\tilde{\mathbf{B}}^{i-2-j}],$$

822 where \otimes denotes the Kronecker product.

823 *Proof.* Firstly, the derivative of matrix power can be obtained using the following equation (Magnus and Neudecker, 2019),

$$824 \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} \tilde{\mathbf{B}}^k = k \mathbf{K}_{dd} \sum_{j=0}^{k-2} [\tilde{\mathbf{B}}^j]^\top \otimes [\tilde{\mathbf{B}}^{k-2-j}], \quad (36)$$

825 where \otimes denotes the Kronecker product. Thus the Hessian of analytic DAG constraints can be obtained as follows:

$$826 \begin{aligned} \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f(\tilde{\mathbf{B}}) &= \sum_{i=0}^{\infty} c_i \text{tr} \nabla_{\tilde{\mathbf{B}}}^2 \tilde{\mathbf{B}}^i \\ &= \mathbf{K}_{dd} \sum_{i=2}^{\infty} i c_i \sum_{j=0}^{i-2} [\tilde{\mathbf{B}}^j]^\top \otimes [\tilde{\mathbf{B}}^{i-2-j}]. \end{aligned} \quad (37)$$

827 □

828 A.8 PROOF OF PROPOSITION 7

829 **Proposition 7.** *For two analytic function $f_1(x) = c_{0,1} + \sum_{i=1}^{\infty} c_{i,1} x^i \in \mathcal{F}$ and $f_2(x) = c_{0,2} + \sum_{i=1}^{\infty} c_{i,2} x^i \in \mathcal{F}$, if $\forall i \geq 1$ we have $c_{i,1} \geq c_{i,2}$, then*

$$830 \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_1(\tilde{\mathbf{B}})) \geq \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_2(\tilde{\mathbf{B}})),$$

831 where $\rho(\cdot)$ denotes the spectral radius of a matrix.

864 *Proof.* Obviously, each entries in the Hessian of $\text{tr} f_1(\tilde{\mathbf{B}})$ is larger than the corresponding ones in
 865 $\text{tr} f_2(\tilde{\mathbf{B}})$. Thus for any unit length vector \mathbf{u} with all positive entries we would have
 866

$$867 \quad \mathbf{u}^\top \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_1(\tilde{\mathbf{B}}) \mathbf{u} \geq \mathbf{u}^\top \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_2(\tilde{\mathbf{B}}) \mathbf{u},$$

868 and then it would be straightforward that
 869

$$870 \quad \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_1(\tilde{\mathbf{B}})) = \max_{\mathbf{u}: \mathbf{u} \geq 0, \|\mathbf{u}\|_2=1} \mathbf{u}^\top \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_1(\tilde{\mathbf{B}}) \mathbf{u}$$

$$871 \quad \geq \max_{\mathbf{u}: \mathbf{u} \geq 0, \|\mathbf{u}\|_2=1} \mathbf{u}^\top \nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_2(\tilde{\mathbf{B}}) \mathbf{u} = \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr} f_2(\tilde{\mathbf{B}}))$$

874 \square

876 B HYPER PARAMETERS

878 In terms of hyper-parameters, our selection involves $\alpha = 0.1$, $\lambda_1 = 0.1$, and $T = 5$. For s we use the
 879 same annealing approach as Bello et al. (2022), but with our strategy to reset s when candidate graph
 880 goes out of the desired region.
 881

882 In all experiments in this paper, for continuous based approaches we use exactly the same hyper
 883 parameter as Bello et al. (2022), for conditional independent test and score based approaches we use
 884 the default parameter in Causal Discovery Toolbox³.

886 C EXTRA EXPERIMENTAL RESULTS

888 In this section, we provide additional experimental results, including true positive rate, false detection
 889 rate and running time for large scale graphs, as well as experimental results on small scale graphs.
 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 ³<https://fentechsolutions.github.io/CausalDiscoveryToolbox/html/index.html>

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Table 5: DAG learning performance (measured in true positive rate, the higher the better, best results in **bold**) of different algorithms on large scale (500-2000 nodes) graphs with different noise distributions. Our algorithm performs better than previous approaches.

Graphs	Nodes	DAGMA	Order-1	Order-2	Order-3	Order-4
ER2-gauss	500	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
ER2-gauss	1000	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00
ER2-gauss	2000	0.96 ± 0.01	0.97 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00
ER3-gauss	500	0.95 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
ER3-gauss	1000	0.94 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
ER3-gauss	2000	0.94 ± 0.01	0.95 ± 0.01	0.95 ± 0.01	0.96 ± 0.01	0.96 ± 0.00
ER4-gauss	500	0.92 ± 0.02	0.94 ± 0.02	0.94 ± 0.02	0.94 ± 0.02	0.95 ± 0.02
ER4-gauss	1000	0.89 ± 0.02	0.92 ± 0.01	0.92 ± 0.01	0.93 ± 0.01	0.93 ± 0.01
ER4-gauss	2000	0.89 ± 0.01	0.90 ± 0.01	0.90 ± 0.01	0.91 ± 0.01	0.91 ± 0.01
ER2-exp	500	0.96 ± 0.01	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
ER2-exp	1000	0.97 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00
ER2-exp	2000	0.00 ± 0.00	0.58 ± 0.48	0.10 ± 0.29	0.10 ± 0.29	0.10 ± 0.29
ER3-exp	500	0.95 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
ER3-exp	1000	0.94 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
ER3-exp	2000	0.00 ± 0.00	0.09 ± 0.28	0.09 ± 0.28	0.10 ± 0.29	0.00 ± 0.00
ER4-exp	500	0.91 ± 0.02	0.93 ± 0.02	0.93 ± 0.02	0.94 ± 0.02	0.95 ± 0.01
ER4-exp	1000	0.89 ± 0.01	0.92 ± 0.01	0.92 ± 0.01	0.93 ± 0.01	0.94 ± 0.01
ER4-exp	2000	0.00 ± 0.00	0.09 ± 0.27	0.09 ± 0.27	0.09 ± 0.27	0.00 ± 0.00
ER2-gumbel	500	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.00
ER2-gumbel	1000	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.00	0.99 ± 0.00	0.99 ± 0.00
ER2-gumbel	2000	0.00 ± 0.00	0.59 ± 0.48	0.10 ± 0.29	0.10 ± 0.30	0.10 ± 0.30
ER3-gumbel	500	0.95 ± 0.01	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
ER3-gumbel	1000	0.95 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.98 ± 0.01
ER3-gumbel	2000	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.10 ± 0.29
ER4-gumbel	500	0.93 ± 0.02	0.95 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	0.96 ± 0.02
ER4-gumbel	1000	0.90 ± 0.01	0.94 ± 0.01	0.94 ± 0.01	0.95 ± 0.01	0.95 ± 0.01
ER4-gumbel	2000	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.09 ± 0.28
SF2-gauss	500	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF2-gauss	1000	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00
SF2-gauss	2000	0.96 ± 0.01	0.97 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00
SF3-gauss	500	0.95 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
SF3-gauss	1000	0.94 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
SF3-gauss	2000	0.94 ± 0.01	0.95 ± 0.01	0.95 ± 0.01	0.96 ± 0.01	0.96 ± 0.00
SF4-gauss	500	0.92 ± 0.02	0.94 ± 0.02	0.94 ± 0.02	0.94 ± 0.02	0.95 ± 0.02
SF4-gauss	1000	0.89 ± 0.02	0.92 ± 0.01	0.92 ± 0.01	0.93 ± 0.01	0.93 ± 0.01
SF4-gauss	2000	0.89 ± 0.01	0.90 ± 0.01	0.90 ± 0.01	0.91 ± 0.01	0.91 ± 0.01
SF2-exp	500	0.96 ± 0.01	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF2-exp	1000	0.97 ± 0.01	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00	0.98 ± 0.00
SF3-exp	500	0.95 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
SF3-exp	1000	0.94 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.97 ± 0.01
SF4-exp	500	0.91 ± 0.02	0.93 ± 0.02	0.93 ± 0.02	0.94 ± 0.02	0.95 ± 0.01
SF4-exp	1000	0.89 ± 0.01	0.92 ± 0.01	0.92 ± 0.01	0.93 ± 0.01	0.94 ± 0.01
SF2-gumbel	500	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.00
SF2-gumbel	1000	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.00	0.99 ± 0.00	0.99 ± 0.00
SF3-gumbel	500	0.95 ± 0.01	0.97 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF3-gumbel	1000	0.95 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.97 ± 0.01	0.98 ± 0.01
SF4-gumbel	500	0.93 ± 0.02	0.95 ± 0.01	0.96 ± 0.01	0.96 ± 0.01	0.96 ± 0.02
SF4-gumbel	1000	0.90 ± 0.01	0.94 ± 0.01	0.94 ± 0.01	0.95 ± 0.01	0.95 ± 0.01

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table 6: DAG learning performance (measured in false detection rate, the lower the better, best results in **bold**) of different algorithms on large scale (500-2000 nodes) graphs with different noise distributions. Our algorithm performs better than previous approaches.

Graphs	Nodes	DAGMA	Order-1	Order-2	Order-3	Order-4
ER2-gauss	500	0.02 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
ER2-gauss	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
ER2-gauss	2000	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.00
ER3-gauss	500	0.04 ± 0.02	0.04 ± 0.02	0.03 ± 0.02	0.03 ± 0.01	0.04 ± 0.02
ER3-gauss	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
ER3-gauss	2000	0.06 ± 0.01	0.06 ± 0.02	0.05 ± 0.01	0.04 ± 0.01	0.04 ± 0.01
ER4-gauss	500	0.08 ± 0.03	0.08 ± 0.03	0.08 ± 0.04	0.07 ± 0.04	0.07 ± 0.03
ER4-gauss	1000	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.02	0.09 ± 0.03	0.10 ± 0.03
ER4-gauss	2000	0.14 ± 0.01	0.13 ± 0.01	0.13 ± 0.02	0.12 ± 0.02	0.12 ± 0.01
ER2-exp	500	0.03 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.02	0.02 ± 0.02
ER2-exp	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
ER3-exp	500	0.05 ± 0.02	0.04 ± 0.02	0.03 ± 0.01	0.04 ± 0.02	0.04 ± 0.02
ER3-exp	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
ER4-exp	500	0.09 ± 0.04	0.09 ± 0.04	0.09 ± 0.05	0.07 ± 0.04	0.06 ± 0.03
ER4-exp	1000	0.12 ± 0.03	0.11 ± 0.02	0.11 ± 0.02	0.09 ± 0.03	0.10 ± 0.02
ER2-gumbel	500	0.03 ± 0.02	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
ER2-gumbel	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.00	0.01 ± 0.01	0.02 ± 0.01
ER3-gumbel	500	0.06 ± 0.02	0.04 ± 0.02	0.03 ± 0.03	0.03 ± 0.01	0.04 ± 0.02
ER3-gumbel	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.02	0.04 ± 0.01	0.03 ± 0.01
ER4-gumbel	500	0.10 ± 0.04	0.09 ± 0.05	0.09 ± 0.04	0.08 ± 0.04	0.08 ± 0.05
ER4-gumbel	1000	0.14 ± 0.03	0.12 ± 0.03	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.03
SF2-gauss	500	0.02 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
SF2-gauss	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
SF2-gauss	2000	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.00
SF3-gauss	500	0.04 ± 0.02	0.04 ± 0.02	0.03 ± 0.02	0.03 ± 0.01	0.04 ± 0.02
SF3-gauss	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
SF3-gauss	2000	0.06 ± 0.01	0.06 ± 0.02	0.05 ± 0.01	0.04 ± 0.01	0.04 ± 0.01
SF4-gauss	500	0.08 ± 0.03	0.08 ± 0.03	0.08 ± 0.04	0.07 ± 0.04	0.07 ± 0.03
SF4-gauss	1000	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.02	0.09 ± 0.03	0.10 ± 0.03
SF4-gauss	2000	0.14 ± 0.01	0.13 ± 0.01	0.13 ± 0.02	0.12 ± 0.02	0.12 ± 0.01
SF2-exp	500	0.03 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.02	0.02 ± 0.02
SF2-exp	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
SF3-exp	500	0.05 ± 0.02	0.04 ± 0.02	0.03 ± 0.01	0.04 ± 0.02	0.04 ± 0.02
SF3-exp	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
SF4-exp	500	0.09 ± 0.04	0.09 ± 0.04	0.09 ± 0.05	0.07 ± 0.04	0.06 ± 0.03
SF4-exp	1000	0.12 ± 0.03	0.11 ± 0.02	0.11 ± 0.02	0.09 ± 0.03	0.10 ± 0.02
SF2-gumbel	500	0.03 ± 0.02	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
SF2-gumbel	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.00	0.01 ± 0.01	0.02 ± 0.01
SF3-gumbel	500	0.06 ± 0.02	0.04 ± 0.02	0.03 ± 0.03	0.03 ± 0.01	0.04 ± 0.02
SF3-gumbel	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.02	0.04 ± 0.01	0.03 ± 0.01
SF4-gumbel	500	0.10 ± 0.04	0.09 ± 0.05	0.09 ± 0.04	0.08 ± 0.04	0.08 ± 0.05
SF4-gumbel	1000	0.14 ± 0.03	0.12 ± 0.03	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.03

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Table 7: DAG learning performance (measured in running time (seconds), the lower the better, best results in **bold**) of different algorithms on large scale (500-2000 nodes) graphs with different noise distributions. Our algorithm performs better than previous approaches.

Graphs	Nodes	DAGMA	Order-1	Order-2	Order-3	Order-4
ER2-gauss	500	171.79 ± 18.30	294.19 ± 62.47	493.73 ± 28.79	482.01 ± 17.59	512.86 ± 33.87
ER2-gauss	1000	364.95 ± 38.47	687.38 ± 69.45	1261.33 ± 47.66	1256.89 ± 33.76	1329.92 ± 315.91
ER2-gauss	2000	1187.17 ± 108.06	3515.56 ± 310.54	6063.65 ± 554.57	6300.09 ± 197.69	6360.65 ± 321.65
ER3-gauss	500	238.54 ± 58.34	394.95 ± 110.87	516.49 ± 14.45	524.93 ± 48.79	528.08 ± 49.59
ER3-gauss	1000	501.24 ± 59.30	1004.85 ± 162.25	1365.24 ± 21.10	1355.23 ± 139.98	1349.25 ± 84.01
ER3-gauss	2000	1636.23 ± 101.64	4903.56 ± 782.06	6037.67 ± 765.47	7072.04 ± 2168.01	6969.54 ± 724.52
ER4-gauss	500	347.06 ± 55.40	519.77 ± 56.88	532.34 ± 10.16	517.46 ± 12.95	534.98 ± 52.22
ER4-gauss	1000	798.01 ± 27.84	1328.32 ± 48.32	1318.45 ± 62.71	1348.08 ± 15.79	1312.53 ± 138.14
ER4-gauss	2000	2461.67 ± 1.49	6520.96 ± 311.23	6560.10 ± 13.55	7666.47 ± 2151.48	7067.87 ± 1011.11
ER2-exp	500	167.08 ± 25.75	291.87 ± 63.66	492.14 ± 30.49	496.83 ± 20.55	504.17 ± 22.59
ER2-exp	1000	359.97 ± 28.62	708.23 ± 61.00	1245.71 ± 71.43	1279.88 ± 41.63	1277.89 ± 45.39
ER3-exp	500	235.71 ± 56.44	440.70 ± 189.91	564.88 ± 145.09	559.54 ± 144.42	550.44 ± 129.03
ER3-exp	1000	515.88 ± 83.18	1100.18 ± 226.98	1371.65 ± 136.38	1401.29 ± 283.39	1503.85 ± 321.35
ER4-exp	500	358.94 ± 46.65	510.53 ± 46.91	513.70 ± 49.38	522.16 ± 15.09	508.01 ± 34.87
ER4-exp	1000	778.40 ± 51.69	1344.07 ± 26.71	1324.00 ± 111.40	1347.33 ± 21.06	1298.91 ± 100.30
ER2-gumbel	500	161.36 ± 24.90	255.86 ± 33.56	501.55 ± 85.14	490.99 ± 11.06	501.44 ± 23.26
ER2-gumbel	1000	330.53 ± 39.10	656.98 ± 62.61	1245.11 ± 47.50	1276.58 ± 15.41	1266.35 ± 38.53
ER3-gumbel	500	232.45 ± 50.71	381.03 ± 85.67	521.98 ± 8.54	514.79 ± 12.92	506.95 ± 18.01
ER3-gumbel	1000	525.92 ± 93.20	1013.67 ± 168.84	1331.88 ± 99.25	1302.93 ± 101.85	1369.44 ± 39.47
ER4-gumbel	500	366.06 ± 33.40	514.95 ± 57.57	540.93 ± 17.35	530.18 ± 20.49	519.87 ± 13.98
ER4-gumbel	1000	805.91 ± 31.96	1367.02 ± 36.36	1260.93 ± 137.27	1434.41 ± 154.30	1335.20 ± 61.78
SF2-gauss	500	171.79 ± 18.30	294.19 ± 62.47	493.73 ± 28.79	482.01 ± 17.59	512.86 ± 33.87
SF2-gauss	1000	364.95 ± 38.47	687.38 ± 69.45	1261.33 ± 47.66	1256.89 ± 33.76	1329.92 ± 315.91
SF2-gauss	2000	1187.17 ± 108.06	3515.56 ± 310.54	6063.65 ± 554.57	6300.09 ± 197.69	6360.65 ± 321.65
SF3-gauss	500	238.54 ± 58.34	394.95 ± 110.87	516.49 ± 14.45	524.93 ± 48.79	528.08 ± 49.59
SF3-gauss	1000	501.24 ± 59.30	1004.85 ± 162.25	1365.24 ± 21.10	1355.23 ± 139.98	1349.25 ± 84.01
SF3-gauss	2000	1636.23 ± 101.64	4903.56 ± 782.06	6037.67 ± 765.47	7072.04 ± 2168.01	6969.54 ± 724.52
SF4-gauss	500	347.06 ± 55.40	519.77 ± 56.88	532.34 ± 10.16	517.46 ± 12.95	534.98 ± 52.22
SF4-gauss	1000	798.01 ± 27.84	1328.32 ± 48.32	1318.45 ± 62.71	1348.08 ± 15.79	1312.53 ± 138.14
SF4-gauss	2000	2461.67 ± 1.49	6520.96 ± 311.23	6560.10 ± 13.55	7666.47 ± 2151.48	7067.87 ± 1011.11
SF2-exp	500	167.08 ± 25.75	291.87 ± 63.66	492.14 ± 30.49	496.83 ± 20.55	504.17 ± 22.59
SF2-exp	1000	359.97 ± 28.62	708.23 ± 61.00	1245.71 ± 71.43	1279.88 ± 41.63	1277.89 ± 45.39
SF3-exp	500	235.71 ± 56.44	440.70 ± 189.91	564.88 ± 145.09	559.54 ± 144.42	550.44 ± 129.03
SF3-exp	1000	515.88 ± 83.18	1100.18 ± 226.98	1371.65 ± 136.38	1401.29 ± 283.39	1503.85 ± 321.35
SF4-exp	500	358.94 ± 46.65	510.53 ± 46.91	513.70 ± 49.38	522.16 ± 15.09	508.01 ± 34.87
SF4-exp	1000	778.40 ± 51.69	1344.07 ± 26.71	1324.00 ± 111.40	1347.33 ± 21.06	1298.91 ± 100.30
SF2-gumbel	500	161.36 ± 24.90	255.86 ± 33.56	501.55 ± 85.14	490.99 ± 11.06	501.44 ± 23.26
SF2-gumbel	1000	330.53 ± 39.10	656.98 ± 62.61	1245.11 ± 47.50	1276.58 ± 15.41	1266.35 ± 38.53
SF3-gumbel	500	232.45 ± 50.71	381.03 ± 85.67	521.98 ± 8.54	514.79 ± 12.92	506.95 ± 18.01
SF3-gumbel	1000	525.92 ± 93.20	1013.67 ± 168.84	1331.88 ± 99.25	1302.93 ± 101.85	1369.44 ± 39.47
SF4-gumbel	500	366.06 ± 33.40	514.95 ± 57.57	540.93 ± 17.35	530.18 ± 20.49	519.87 ± 13.98
SF4-gumbel	1000	805.91 ± 31.96	1367.02 ± 36.36	1260.93 ± 137.27	1434.41 ± 154.30	1335.20 ± 61.78

	PC	GES	DAGMA	Exponential	Order 1	Order 2	Order 3	Order 4
SHD	563.9 ± 23.84	4490.2 ± 62.52	588.8 ± 18.33	488.6 ± 24.29	429.6 ± 24.73	410.6 ± 15.25	401.0 ± 16.64	389.4 ± 16.70
				Exp MLE	Order 1 MLE	Order 2 MLE	Order 3 MLE	Order 4 MLE
SHD				518.00 ± 23.02	453.70 ± 42.12	447.30 ± 51.85	409.50 ± 31.02	433.00 ± 68.98
				PC Exp	PC Order-1	PC Order-2	PC Order-3	PC Order-4
SHD				275.40 ± 16.01	274.40 ± 15.44	273.10 ± 15.72	271.80 ± 14.75	276.00 ± 14.66
				PC Exp MLE	PC Order-1 MLE	PC Order-2 MLE	PC Order-3 MLE	PC Order-4 MLE
SHD				274.30 ± 14.71	284.30 ± 19.43	272.20 ± 14.04	273.00 ± 17.79	270.20 ± 12.58

	PC	GES	DAGMA	Exponential	Order 1	Order 2	Order 3	Order 4
SHDC	321.30 ± 27.77	4626.20 ± 69.05	674.00 ± 31.09	588.60 ± 59.81	466.50 ± 26.43	458.40 ± 30.85	447.20 ± 30.81	439.90 ± 37.06
				Exp MLE	Order 1 MLE	Order 2 MLE	Order 3 MLE	Order 4 MLE
SHDC				574.50 ± 42.84	490.80 ± 66.99	486.80 ± 76.47	444.30 ± 42.22	479.50 ± 100.71
				PC Exp	PC Order-1	PC Order-2	PC Order-3	PC Order-4
SHDC				236.20 ± 15.16	236.20 ± 15.16	236.20 ± 15.16	236.20 ± 15.16	236.20 ± 15.16
				PC Exp MLE	PC Order-1 MLE	PC Order-2 MLE	PC Order-3 MLE	PC Order-4 MLE
SHDC				231.30 ± 13.64	257.50 ± 26.14	236.10 ± 16.23	236.80 ± 23.77	231.60 ± 13.17

Table 11: DAG learning performance (measured in structural hamming distance, the lower the better, best results in **bold**) of different algorithms on 1000-node ER1 graphs with Gaussian noise with observation data normalized. Our algorithms performs better than the previous approaches, and as higher order DAG constraints suffers less to gradient vanishing, it tends to have better performance. We compare differential DAG learning approaches with conditional independent test based PC (Spirtes and Glymour, 1991) algorithm and score based GES (Chickering, 2002) algorithm. The result is reported in the format of average ± standard derivation gathered from 10 different simulations. The results are reported as averages ± standard deviations, calculated from 10 independent simulations. In addition to the MSE score function, we also applied the MLE score function described in Ng et al. (2020). Furthermore, rather than only considering edges between variables with correlation coefficients greater than 0.1, we also evaluated cases where edges are restricted to those in the PC-estimated CPDAG (algorithms whose names begin with 'PC').

D IMPLEMENTATION FOR ALGORITHM 1

```

def _h_grad(self, W, s, eps=1e-20):
    M_ = W * W / s
    Iw = self.Id - M_ # self.Id is identity matrix
    icnt = 1
    Inv = self.Id + M_
    while icnt < 2 * self.d:
        M_ = M_ @ M_
        Inv = Inv + Inv @ M_
        icnt *= 2
        if self.np.max(self.np.abs(M_)) < eps:
            break
        if self.np.any(self.np.isnan(Inv)):
            break

    if self.np.any(self.np.isinf(Inv)):
        return self.np.zeros_like(Inv)

    if self.np.any(self.np.isnan(Inv)):
        return self.np.zeros_like(Inv)

    return Inv / s

def compute_h_grad(self, W, s):

```


1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

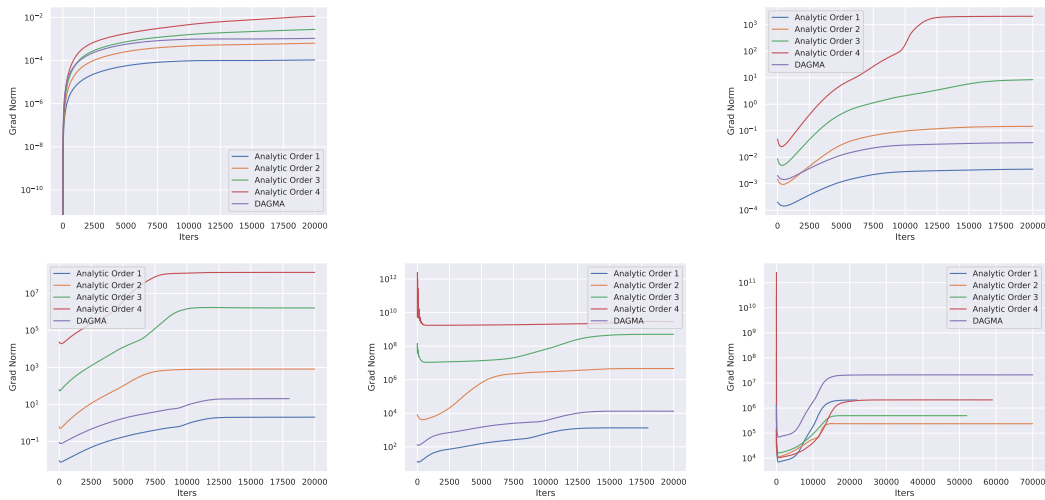


Figure 1: Gradient Norm v.s. Optimization Iterations. **Top:** Frobenius norm of gradients vs. gradient descent steps for the first two iterations in Algorithm 2. **Bottom:** Frobenius norm of gradients vs. gradient descent steps for the last three iterations in Algorithm 2. In most cases, our higher-order DAG constraints exhibit larger gradient norms compared to DAGMA, enabling our algorithm to often converge to better solutions than DAGMA.

```

M = self._h_grad(W, s)
if self.np.any(self.np.isnan(M)) or self.np.linalg.norm(
    M @ (s * self.Id - W * W) - self.Id, ord='fro') >
    1e-6:
    if isinstance(W, cupy.ndarray):
        _, s, v = cupy.linalg.svd(W * W) # cupy does not
        have a eig lib, thus use spectral norm as an
        estimation
        cs = cupy.max(s) + 0.1 * self.h_order
    else:
        cs = np.max(np.abs(np.linalg.eigvals(
            W * W))) + 0.1 * self.h_order
else:
    cs = s
return M, cs

```

E EMPIRICAL RESULTS ON GRADIENTS VANISHING

In this section, we present empirical results on gradient issues in DAG learning. According to Proposition 5, the gradients of higher-order DAG constraints should have larger norms than those of lower-order constraints, given the same candidate adjacency matrix. Additionally, the behavior of our Order-1 DAG constraints is expected to align closely with that of DAGMA.

In the path-following algorithm described in Algorithm 2, five iterations are used to tune the scale of the score function. During each iteration, tens of thousands of gradient descent steps are performed. We recorded the gradient norms for various DAG constraints over the five iterations using a 1000-node ER3 DAG learning problem with Gaussian noise. The results are shown in Figure 1. In most cases, our higher-order DAG constraints exhibit larger gradient norms compared to DAGMA, enabling our algorithm to often converge to better solutions than DAGMA.