
ANALYTIC DAG CONSTRAINTS FOR DIFFERENTIABLE DAG LEARNING

005 **Anonymous authors**

006 Paper under double-blind review

ABSTRACT

011 Recovering underlying Directed Acyclic Graph (DAG) structures from observational
012 data presents a formidable challenge due to the combinatorial nature of
013 the DAG-constrained optimization problem. Recently, researchers have identified
014 gradient vanishing as one of the primary obstacles in differentiable DAG
015 learning and have proposed several DAG constraints to mitigate this issue. By
016 developing the necessary theory to establish a connection between analytic functions
017 and DAG constraints, we demonstrate that analytic functions from the set
018 $\{f(x) = c_0 + \sum_{i=1}^r c_i x^i \mid \forall i > 0, c_i > 0; r = \lim_{i \rightarrow \infty} c_i/c_{i+1} > 0\}$ can be
019 employed to formulate effective DAG constraints. Furthermore, we establish that
020 this set of functions is closed under several functional operators, including differ-
021 entiation, summation, and multiplication. Consequently, these operators can be
022 leveraged to create novel DAG constraints based on existing ones. Using these
023 properties, we designed a series of DAG constraints and developed an efficient
024 algorithm to evaluate these DAG constraints. Experiments conducted in various
025 settings demonstrate that our DAG constraints outperform previous state-of-the-art
026 approaches.

1 INTRODUCTION

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

034 There are two main categories of DAG learning approaches: constraint-based and score-based
035 methods. Most constraint-based approaches, e.g., PC (Spirtes and Glymour, 1991), FCI (Spirtes et al.,
036 1995; Colombo et al., 2012), rely on conditional independence tests, which typically necessitate a
037 large sample size (Shah and Peters, 2020; Vowels et al., 2021). The score-based approaches, including
038 exact methods based on dynamic programming (Koivisto and Sood, 2004; Singh and Moore, 2005;
039 Silander and Myllymäki, 2006), A* search (Yuan et al., 2011; Yuan and Malone, 2013), and integer
040 programming (Cussens, 2011), as well as greedy methods like GES (Chickering, 2002), model
041 the validity of a graph according to some score function and are often formulated and solved as
042 discrete optimization problems. A key challenge for score-based methods is the super-exponential
043 combinatorial search space of DAGs w.r.t number of nodes (Chickering, 1996; Chickering et al.,
044 2004).

045 Recently, Zheng et al. (2018) developed a continuous DAG learning approach using Langrange
046 Multiplier methods and a differentiable DAG constraint based on the trace of the matrix exponential
047 of the weighted adjacency matrix. The resulting method, named NOTEARS, demonstrated superior
048 performance in estimating linear DAGs with equal noise variances. Very recently, Zhang et al. (2022)
049 and Bello et al. (2022) suggest that one main issue for NOTEARS and its derivatives, such as Yu et al.
050 (2019), is gradient vanishing for linear DAG models with equal variance. They have thus proposed
051 new continuous DAG constraints by based on geometric series of matrices as well as log-determinant
052 of matrices.

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

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

061 The answer is yes. We demonstrate that any analytic function within the class of functions denoted
 062 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
 063 formulate Directed Acyclic Graph (DAG) constraints. In fact, the DAG constraints introduced in
 064 Zheng et al. (2018), Zhang et al. (2022), and Bello et al. (2022) can all be interpreted as being
 065 based on analytic functions from \mathcal{F} . Furthermore, we establish that the function class \mathcal{F} remains
 066 closed under various function operators, including differentiation, function addition, and function
 067 multiplication. Leveraging this insight, we can construct novel DAG constraints based on pre-
 068 existing ones. Additionally, we can analyze the performance of these derived DAG constraints using
 069 techniques rooted in analytic functions.

071 2 PRELIMINARIES

072 **DAG and Linear SEM** Given a directed acyclic graph (DAG) \mathcal{G} defined over random vector $\mathbf{x} =$
 073 $[x_1, x_2, \dots, x_d]^{\top}$, the corresponding distribution $P(\mathbf{x})$ is assumed to satisfy the Markov assumption
 074 (Spirtes et al., 2000; Pearl, 2000). We consider \mathbf{x} to follow a linear Structural Equation Model (SEM):
 075

$$076 \mathbf{x} = \mathbf{B}^{\top} \mathbf{x} + \mathbf{e}. \quad (1)$$

077 Here, $\mathbf{B} \in \mathbb{R}^{d \times d}$ represents the weighted adjacency matrix that characterizes the DAG \mathcal{G} , and
 078 $\mathbf{e} = [e_1, e_2, \dots, e_d]^{\top}$ represents the exogenous noise vector, comprising d independent random
 079 variables. To simplify notation, we use $\mathcal{G}(\mathbf{B})$ to denote the graph induced by the weighted adjacency
 080 matrix \mathbf{B} , and we interchangeably use the terms ‘random variables’ and ‘vertices’ or ‘nodes’.
 081

082 We aim to estimate the DAG \mathcal{G} from n i.i.d. observational examples of \mathbf{x} , denoted by $\mathbf{X} \in \mathbb{R}^{n \times d}$.
 083 Generally, the DAG \mathcal{G} can be identified only up to its Markov equivalence class under the faithfulness
 084 (Spirtes et al., 2000) or the sparsest Markov representation assumption (Raskutti and Uhler, 2018). It
 085 has been demonstrated that for linear SEMs with homoscedastic errors, where the noise terms are
 086 specified up to a constant (Loh and Bühlmann, 2013), and for linear non-Gaussian SEMs, where
 087 no more than one of the noise terms is Gaussian (Shimizu et al., 2006), the true DAG can be fully
 088 identified. In our study, we specifically focus on linear SEMs with equal noise variances (Peters and
 089 Bühlmann, 2013), where the scale of the data may be either known or unknown. When the scale is
 090 known, it is possible to fully recover the DAG. However, in the case of an unknown scale, the DAG
 091 may only be identified up to its Markov equivalence class.
 092

093 **Continuous DAG learning** In recent years, a series of continuous Directed Acyclic Graph (DAG)
 094 learning algorithms Bello et al. (2022); Ng et al. (2020); Zhang et al. (2022); Yu et al. (2021; 2019);
 095 Zheng et al. (2018) has been introduced, demonstrating superior performance when applied to linear
 096 Structural Equation Models (SEMs) with equal noise variances and known data scale. These methods
 097 can be expressed as follows:
 098

$$100 \underset{\mathbf{B}}{\operatorname{argmin}} S(\mathbf{B}, \mathbf{X}), \text{ s.t. } h(\mathbf{B}) = 0. \quad (2)$$

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

108 3 ANALYTIC DAG CONSTRAINTS 109

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

115 3.1 ANALYTIC FUNCTIONS AS DAG CONSTRAINTS 116

117 In mathematics, a power series

$$118 \quad f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i, \quad (3)$$

120 which converges for $|x| < r = \lim_{i \rightarrow \infty} |c_i/c_{i+1}|$, defines an analytic function f on the open interval
121 $(-r, r)$, and r is known as the convergence radius. When we replace x with a square matrix \mathbf{A} , we
122 obtain an analytic function f of a matrix as follows:

$$124 \quad f(\mathbf{A}) = c_0 \mathbf{I} + \sum_{i=1}^{\infty} c_i \mathbf{A}^i, \quad (4)$$

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

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

$$131 \quad \mathcal{F} = \{f | f(x) = c_0 + \sum_{i=1}^{\infty} c_i x^i; \forall i > 0, c_i > 0; \lim_{i \rightarrow \infty} c_i/c_{i+1} > 0\}, \quad (5)$$

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

135 **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
136 \mathcal{G} , and let f be an analytic function in the form of (4), where we further assume $\forall i > 0$ we have
137 $c_i > 0$, then \mathcal{G} is acyclic if and only if

$$139 \quad \text{tr}[f(\tilde{\mathbf{B}})] = c_0 d. \quad (6)$$

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

144 **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
145 of $\text{tr}[f(\tilde{\mathbf{B}})]$ w.r.t. $\tilde{\mathbf{B}}$ is

$$147 \quad \nabla_{\tilde{\mathbf{B}}} \text{tr}[f(\tilde{\mathbf{B}})] = [\nabla_x f(x)|_{x=\tilde{\mathbf{B}}}]^\top. \quad (7)$$

149 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
150 and necessary to enforce DAGness (Wei et al., 2020; Ng et al., 2022). Meanwhile, evaluating matrix
151 polynomials efficiently is highly nontrivial (Higham, 2008). For matrix analytic functions such as
152 exponentials or logarithms, however, efficient algorithms exist (Higham, 2008).

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

$$160 \quad \text{Constraints: } \text{tr}[\exp(\tilde{\mathbf{B}})] = \sum_{i=0}^{\infty} \tilde{\mathbf{B}}^i / i! = d, \quad \text{Gradient: } \nabla_{\tilde{\mathbf{B}}} \exp(\tilde{\mathbf{B}}) = \exp(\tilde{\mathbf{B}})^\top, \quad (8)$$

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

Recently Bello et al. (2022) and Zhang et al. (2022) have suggested that exponential-based DAG constraints suffers from gradient vanishing. One cause of gradient vanishing arises from the small coefficients of high order terms. The convergence radius for the exponential is ∞ , that 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, the higher order terms contribute almost nothing in the DAG constraints, which indicates that it would not be efficient to prohibit possible long loops in candidate adjacency matrices.

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

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

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

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

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

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

3.2 CONSTRUCTING DAG CONSTRAINTS BY FUNCTIONAL OPERATOR

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

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

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

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 larger than 1, then $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with spectral radius $\rho(\hat{\mathbf{B}}) \leq r$ forms a DAG if and only if

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

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

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-jancency matrix $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with spectral radius $\rho(\tilde{\mathbf{B}}) \leq \min(\lim_{i \rightarrow \infty} c_i^1/c_{i+1}^1, \lim_{i \rightarrow \infty} c_i^2/c_{i+1}^2)\}$, the following three statements are equivalent:

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

216 3. $\text{tr}[f_1(\tilde{\mathbf{B}})f_2(\tilde{\mathbf{B}})] = c_0^1c_0^2d.$
 217

218 Particularly for $f_{\log}^s(x)$ and $f_{inv}^s(x)$, due to the specific property of $f_{inv}^s(x)$, we have
 219

220

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

223 which implies that the n^{th} derivative of function $1/(s-x)$ is propositional to the the order- $(n+1)$
 224 power of $1/(s-x)$. Using this property, the value of $(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-1}$ can be cached and then used
 225 to generate a series of DAG constraints as well as their gradients. Similarly, the value of matrix
 226 exponential $\exp((\tilde{\mathbf{B}})/s)$ can also be cached during the evaluation of DAG constraints to accelerate
 227 the computation. Furthermore, the gradients of the DAG constraints will also increase as n increases.
 228

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

231 $\text{tr}[(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n}] = d.$

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

234

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

236 where k is an arbitrary positive integer, and $\|\cdot\|$ denote an arbitrary matrix norm induced by vector
 237 p -norm.
 238

239 **Gradient Vanishing and Numeric Stability** For the series of DAG constraints constructed from
 240 Equation (13), as gradient vanishing is one of the main challenges for differentiable DAG learning,
 241 according to Proposition 5 we may prefer larger n to achieve better performance in practice.
 242 Furthermore, choosing a smaller s may also help to amplify the gradient of DAG constraints. Therefore,
 243 Bello et al. (2022) applied an annealing strategy on s to improve performance, while Zhang
 244 et al. (2022) used a fixed $s = 1.0$ in their implementation. However, in practice, especially when
 245 incorporating the DAG constraints with first-order optimizers, the spectral radius of the candidate
 246 $\tilde{\mathbf{B}}$ can often be larger than s . Bello et al. (2022) applied a simple heuristics to search for the proper
 247 s , while Zhang et al. (2022) truncated the power series to avoid numerical issues in higher-order
 248 terms. However, in practice, we observed that Zhang et al. (2022) encountered some numerical issues
 249 for large graphs, and the simple heuristics used by Bello et al. (2022) may result in a sacrifice in
 250 performance. Based on our analysis, if $\tilde{\mathbf{B}}$ goes out, it can be verified by checking if the power series
 251 $\sum_{i=0}^{\infty} (\tilde{\mathbf{B}}/s)^i$ converges to $(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-1}$, and s can be chosen based on the spectral radius of $\tilde{\mathbf{B}}$.

252

253 **Algorithm 1** Efficient Evaluation of Gradients

254

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

257

258 1: $\mathbf{D} \leftarrow \mathbf{I} + \tilde{\mathbf{B}}/s, \mathbf{W} \leftarrow \tilde{\mathbf{B}}/s$
 259 2: $k = 1$
 260 3: **while** $\|\mathbf{D}(\mathbf{I} - \tilde{\mathbf{B}}) - \mathbf{I}\| > \epsilon$ and $k < 2d$ **do**
 261 4: $\mathbf{W} \leftarrow \mathbf{W} \times \mathbf{W}$
 262 5: $\mathbf{D} \leftarrow \mathbf{D} \times (\mathbf{W} + \mathbf{I})$
 263 6: $k \leftarrow 2k$
 264 7: **end while**
 265 8: **if** $\|\mathbf{D}(\mathbf{I} - \tilde{\mathbf{B}}) - \mathbf{I}\| > \epsilon$ **then**
 266 9: $s \leftarrow \rho(\tilde{\mathbf{B}}) + \xi$, goto line 1
 267 10: **else**
 268 11: For f_{\log}^s return \mathbf{D}^\top/s
 269 12: For $[f_{inv}^s(\tilde{\mathbf{B}})]^n$ return $n[\mathbf{D}^\top/s]^{n+1}$
 13: **end if**

254

255 **Algorithm 2** Path following algorithm

256

257 **Input:** $\mathbf{X} \in \mathbb{R}^{n \times d}; S; f \in \mathcal{F}; \lambda_1; \mu_0;$
 258 $\alpha \in (0, 1); T_{outer}; T_{inner}; \gamma > 0$

259 **Output:** Estimated \mathbf{B}

260 1: $i \leftarrow 0, \mu \leftarrow \mu_0, \mathbf{B}_0 = \mathbf{0}$
 261 2: **for** $i = 0; i < T_{outer}; i++$ **do**
 262 3: $\mathbf{B}_{i+1} \leftarrow \mathbf{B}_i$ ▷ Optimize over
 $\mu[S(\mathbf{B}, \mathbf{X}) + \lambda_1 \|\mathbf{B}\|_1] + 1/2f(\mathbf{B} \odot \mathbf{B})$
 263 4: **for** $j = 0; j < T_{inner}; j++$ **do**
 264 5: $\tilde{\mathbf{B}} \leftarrow \mathbf{B}_{i+1} \odot \mathbf{B}_{i+1}$
 265 6: $\mathbf{B}_{i+1} \leftarrow \mathbf{B}_{i+1} - \gamma \mu [\nabla_{\mathbf{B}} S(\mathbf{B}, \mathbf{X}) +$
 $\lambda_1 \text{sign}(\mathbf{B})] - \gamma \nabla_{\tilde{\mathbf{B}}} f(\tilde{\mathbf{B}}) \odot \mathbf{B}_{i+1}$
 266 7: **end for**
 267 8: $\mu \leftarrow \mu \times \alpha$
 268 9: $\hat{\mathbf{B}} \leftarrow \mathbf{B}_{i+1}$
 269 10: **end for**
 11: **Return** $\hat{\mathbf{B}}$

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

$$272 \quad 273 \quad 274 \quad \mathbf{L}_t = \sum_{i=0}^t (\tilde{\mathbf{B}}/s)^i, \quad (14)$$

275 then it is evident that

$$276 \quad 277 \quad \mathbf{L}_{2t} = \mathbf{L}_t + (\tilde{\mathbf{B}}/s)^t \mathbf{L}_t, \quad (15)$$

278 which indicates that the term \mathbf{L}_t can be obtained with $\mathcal{O}(\log t)$ time complexity. Furthermore, using
 279 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
 280 strategy for searching s , we can use Algorithm 1 to efficiently compute the DAG constraints.

281 3.3 OVERALL OPTIMIZATION FRAMEWORK

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

$$284 \quad 285 \quad \underset{\mathbf{B}}{\operatorname{argmin}} S(\mathbf{B}, \mathbf{X}), \quad \text{s.t. } \text{tr}f(\mathbf{B} \odot \mathbf{B}) = c_0 d, \rho(\mathbf{B} \odot \mathbf{B}) < r, \quad (16)$$

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

287 In our work, we choose to use the path-following approach with an ℓ_1 regularizer, as in Bello et al.
 288 (2022). This is because in the Lagrange approaches applied in Zhang et al. (2022); Yu et al. (2021);
 289 Zheng et al. (2018); Yu et al. (2019), the Lagrangian multiplier must be set to very large value
 290 to enforce DAGness, which may result in numerical instability. In the path-following approach,
 291 instead of using large Lagrangian multipliers, a small coefficients are added to the score function S
 292 as follows¹

$$293 \quad 294 \quad \underset{\mathbf{B}}{\operatorname{argmin}} \mu[S(\mathbf{B}, \mathbf{X}) + \lambda_1 \|\mathbf{B}\|_1] + \text{tr}f(\mathbf{B} \odot \mathbf{B}), \quad \text{s.t. } \rho(\mathbf{B} \odot \mathbf{B}) < r, \quad (17)$$

295 where λ_1 is the user-specified weight for the ℓ_1 regularizer. For the additional constraints $\rho(\mathbf{B} \odot \mathbf{B}) <$
 296 r , with properly chosen initial value and step-length, it can usually be satisfied. Also it is notable that
 297 $\|\mathbf{B}\|_1 < r$ is a sufficient condition for $\rho(\mathbf{B} \odot \mathbf{B}) < r$, and thus the sparsity constraints also encourage
 298 this condition to be satisfied. Based on Bello et al. (2022), we implemented a path-following shown
 299 in Algorithm 2.

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

307 4 NON-CONVEXITY ANALYSIS OF ANALYTIC DAG CONSTRAINTS

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

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

$$312 \quad 313 \quad 314 \quad \nabla_{\tilde{\mathbf{B}}}^2 \text{tr}f(\tilde{\mathbf{B}}) = \mathbf{K}_{dd} \sum_{i=2}^{\infty} i c_i \sum_{j=0}^{i-2} [\tilde{\mathbf{B}}^j]^{\top} \otimes [\tilde{\mathbf{B}}^{i-2-j}], \quad (18)$$

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

$$317 \quad 318 \quad \mathbf{K}_{d,d} \text{vec}(\mathbf{A}) = \text{vec}(\mathbf{A}^{\top}).$$

319 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
 320 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: $\text{tr} f_{\log}^s(\mathbf{B} \odot \mathbf{B}) = 0$;

Graphs	#Nodes	DAGMA	Order 1	Order 2	Order 3	Order 4
ER2-Gaussian	500	44.90 \pm 32.95	33.40 \pm 23.46	31.70 \pm 19.47	30.60 \pm 19.07	29.80 \pm 20.97
	1000	94.80 \pm 35.80	69.60 \pm 27.64	55.60 \pm 19.13	52.40 \pm 19.86	57.30 \pm 21.39
	2000	235.40 \pm 62.76	176.00 \pm 47.77	153.30 \pm 35.92	135.60 \pm 38.91	131.00 \pm 28.65
ER3-Gaussian	500	125.30 \pm 44.55	101.10 \pm 39.03	90.30 \pm 39.56	93.90 \pm 31.26	92.90 \pm 45.56
	1000	339.60 \pm 67.80	242.80 \pm 72.21	210.30 \pm 60.98	184.90 \pm 47.44	165.80 \pm 35.95
	2000	669.50 \pm 140.61	610.70 \pm 136.84	555.30 \pm 106.01	479.50 \pm 88.72	424.90 \pm 64.39
ER4-Gaussian	500	307.60 \pm 116.53	261.40 \pm 102.81	263.00 \pm 122.34	246.70 \pm 110.28	223.80 \pm 97.46
	1000	878.50 \pm 174.96	689.20 \pm 165.62	695.50 \pm 134.41	619.80 \pm 150.43	626.30 \pm 157.59
	2000	1922.30 \pm 187.69	1785.40 \pm 184.47	1779.30 \pm 211.23	1655.40 \pm 181.75	1574.10 \pm 152.23
ER2-Exp	500	58.20 \pm 31.58	40.50 \pm 26.93	28.90 \pm 16.60	31.00 \pm 25.67	35.20 \pm 34.32
	1000	93.90 \pm 33.96	68.70 \pm 23.20	54.00 \pm 16.26	50.90 \pm 17.29	57.90 \pm 24.93
ER3-Exp	500	142.70 \pm 50.13	106.00 \pm 39.15	95.10 \pm 32.68	99.60 \pm 37.94	100.30 \pm 47.52
	1000	321.10 \pm 83.82	242.80 \pm 68.51	212.40 \pm 67.87	187.60 \pm 61.87	173.10 \pm 49.03
ER4-Exp	500	336.00 \pm 124.19	292.70 \pm 123.41	294.90 \pm 130.66	254.40 \pm 133.05	214.70 \pm 84.29
	1000	879.40 \pm 162.98	718.20 \pm 127.12	710.60 \pm 151.41	640.50 \pm 148.24	619.70 \pm 133.44
ER2-Gumbel	500	45.10 \pm 33.28	22.60 \pm 20.04	21.30 \pm 18.41	19.80 \pm 16.39	16.20 \pm 11.91
	1000	80.50 \pm 42.65	49.90 \pm 24.54	39.90 \pm 14.04	36.80 \pm 15.18	45.90 \pm 23.18
ER3-Gumbel	500	147.10 \pm 54.19	94.10 \pm 40.87	76.60 \pm 60.77	60.60 \pm 31.34	85.30 \pm 50.75
	1000	297.90 \pm 72.40	215.40 \pm 52.35	185.00 \pm 71.98	173.90 \pm 57.09	147.70 \pm 40.51
ER4-Gumbel	500	338.80 \pm 127.56	273.70 \pm 131.13	257.50 \pm 111.06	232.40 \pm 121.98	234.70 \pm 149.59
	1000	919.90 \pm 182.38	722.80 \pm 177.86	734.80 \pm 177.78	620.70 \pm 187.56	564.10 \pm 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 \pm 43.51	24.30 \pm 43.90	32.40 \pm 49.38	34.20 \pm 45.56	41.50 \pm 48.45
	1000	44.90 \pm 34.38	41.20 \pm 36.02	22.50 \pm 13.21	29.20 \pm 20.07	58.10 \pm 27.58
	2000	189.80 \pm 99.47	162.90 \pm 73.30	172.10 \pm 74.35	152.20 \pm 90.29	172.60 \pm 124.12
SF3	500	58.10 \pm 33.90	51.10 \pm 32.10	41.10 \pm 17.91	49.80 \pm 24.58	71.00 \pm 23.46
	1000	169.40 \pm 60.82	158.10 \pm 46.70	161.20 \pm 55.25	162.50 \pm 57.54	195.40 \pm 75.60
	2000	928.70 \pm 148.70	896.00 \pm 101.85	897.10 \pm 146.78	891.50 \pm 143.40	999.70 \pm 206.36
SF4	500	131.20 \pm 42.63	136.80 \pm 41.71	134.40 \pm 39.34	128.90 \pm 36.68	151.60 \pm 37.05
	1000	431.70 \pm 119.22	404.00 \pm 88.89	400.30 \pm 76.49	386.90 \pm 93.16	394.50 \pm 111.57
	2000	1525.10 \pm 299.02	1500.50 \pm 297.88	1444.70 \pm 291.45	1395.60 \pm 264.90	1418.90 \pm 228.86
SF2	500	25.90 \pm 44.45	23.40 \pm 44.41	32.10 \pm 49.08	35.00 \pm 48.84	37.20 \pm 45.21
	1000	43.70 \pm 34.48	41.20 \pm 36.02	32.00 \pm 34.13	29.10 \pm 19.68	59.10 \pm 30.34
SF3	500	57.70 \pm 33.68	57.70 \pm 33.64	41.80 \pm 20.37	43.20 \pm 15.75	66.70 \pm 24.36
	1000	177.10 \pm 67.53	165.40 \pm 57.70	171.60 \pm 66.80	175.10 \pm 69.09	195.90 \pm 80.37
SF4	500	127.50 \pm 40.84	132.80 \pm 40.39	129.90 \pm 42.07	135.50 \pm 44.21	152.40 \pm 39.94
	1000	408.80 \pm 119.71	419.70 \pm 108.01	388.50 \pm 53.01	394.30 \pm 88.95	395.30 \pm 109.37
SF2	500	23.10 \pm 44.78	17.70 \pm 44.51	16.70 \pm 44.15	20.00 \pm 45.75	33.40 \pm 49.30
	1000	29.20 \pm 24.77	24.70 \pm 24.85	12.50 \pm 11.40	16.20 \pm 13.96	47.90 \pm 25.69
SF3	500	33.50 \pm 32.98	25.20 \pm 27.85	19.40 \pm 12.37	19.00 \pm 7.44	50.00 \pm 22.41
	1000	107.50 \pm 50.50	114.50 \pm 59.80	106.60 \pm 64.77	103.70 \pm 58.15	133.60 \pm 88.43
SF4	500	77.70 \pm 41.43	76.20 \pm 41.86	67.90 \pm 26.47	79.20 \pm 23.87	101.40 \pm 22.37
	1000	333.10 \pm 118.06	348.80 \pm 110.93	309.20 \pm 51.86	321.50 \pm 83.13	339.70 \pm 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

435
436 Table 3: DAG learning performance (measured in structural hamming distance, the lower the better, best results
437 in **bold**) of different algorithms on 1000-node ER1 graphs with Gaussian noise with observation data normalized.
438 Our algorithms performs better than the previous approaches, and as higher order DAG constraints suffers less to
439 gradient vanishing, it tends to have better performance. We compare differential DAG learning approaches with
440 conditional independent test based PC (Spirtes and Glymour, 1991) algorithm and score based GES (Chickering,
441 2002) algorithm. The result is reported in the format of average \pm standard derivation gathered from 10 different
442 simulations.

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

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

451 5.1 LINEAR SEM WITH KNOWN GROUND TRUTH SCALE

452 For linear SEM with a known ground truth scale, our experimental setting is similar to Bello et al.
453 (2022); Zhang et al. (2022); Zheng et al. (2018). We generated two different types of random graphs:
454 ER (Erdős-Rényi) and SF (Scale-Free) graphs with different numbers of expected edges. We use ER_n
455 (SF_n) to denote graphs with d nodes and nd expected edges. Edge weights generated from a uniform
456 distribution over the union of two intervals $[-2, -0.5] \cup [0.5, 2.0]$ are assigned to each edge to form
457 a weighted adjacency matrix B . Then, n samples are generated from the linear SEM $x = Bx + e$
458 to form an $n \times d$ data matrix X , where the noise e is iid sampled from Gaussian, Exponential, or
459 Gumbel distribution. As Bello et al. (2022); Zhang et al. (2022); Zheng et al. (2018) achieved nearly
460 perfect results on small and sparse graphs, we considered more challenging large and denser graphs
461 in our experiments. We set the sample size $n = 1000$ and consider 3 different numbers of nodes
462 $d = 500, 1000, 2000$. For each setting, we conducted 10 random simulations to obtain an average
463 performance. All these experiments were performed using an A100 GPU, and all computations were
464 done in double precision. Our algorithms were compared with the previous state-of-the-art approach
465 DAGMA Bello et al. (2022). The original version of DAGMA Bello et al. (2022) used numpy and
466 ran on CPU; we replaced numpy with cupy to get a GPU version of DAGMA, which performed
467 identically to the CPU version.

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

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

482 In terms of running time, all algorithms had similar running times, typically about 5 minutes for a
483 500-node graph, 10-20 minutes for a 1000-node graph, and around 2 hours for a 2000-node graph.
484 Due to limited time and resources, we only considered $d = 2000$ for Gaussian noises, and for other
485 cases, we only considered $d = 500, 1000$.

	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 Nicholas J Higham. *Functions of matrices: theory and computation*. SIAM, 2008.
- 556 Diviyani Kalainathan, Olivier Goudet, and Ritik Dutta. Causal discovery toolbox: Uncovering causal relationships
557 in python. *The Journal of Machine Learning Research*, 21(1):1406–1410, 2020.
- 558
- 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
- 562 Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- 563 Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural
564 DAG learning. In *International Conference on Learning Representations*, 2020.
- 565
- 566 Po Ling Loh and Peter Bühlmann. High-dimensional learning of linear causal networks via inverse covariance
567 estimation. *Journal of Machine Learning Research*, 2013.
- 568
- 569 Jan R Magnus and Heinz Neudecker. The commutation matrix: some properties and applications. *The annals of*
570 *statistics*, 7(2):381–394, 1979.
- 571
- 572 Jan R Magnus and Heinz Neudecker. *Matrix differential calculus with applications in statistics and econometrics*.
573 John Wiley & Sons, 2019.
- 574
- 575 Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and DAG constraints for learning
576 linear DAGs. *Advances in Neural Information Processing Systems*, 33, 2020.
- 577
- 578 Ignavier Ng, Sébastien Lachapelle, Nan Rosemary Ke, Simon Lacoste-Julien, and Kun Zhang. On the conver-
579 gence of continuous constrained optimization for structure learning. In *International Conference on Artificial*
580 *Intelligence and Statistics*, pages 8176–8198. PMLR, 2022.
- 581
- 582 Ignavier Ng, Biwei Huang, and Kun Zhang. Structure learning with continuous optimization: A sober look and
583 beyond. *arXiv preprint arXiv:2304.02146*, 2023.
- 584
- 585 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
586 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep
587 learning library. *Advances in neural information processing systems*, 32, 2019.
- 588
- 589 Judea Pearl. Models, reasoning and inference. *Cambridge, UK: Cambridge University Press*, 19, 2000.
- 590
- 591 Jonas Peters and Peter Bühlmann. Identifiability of Gaussian structural equation models with equal error
592 variances. *Biometrika*, 101(1):219–228, 2013.
- 593
- 594 Garvesh Raskutti and Caroline Uhler. Learning directed acyclic graph models based on sparsest permutations.
595 *Stat*, 7(1):e183, 2018.
- 596
- 597 Alexander Reisach, Christof Seiler, and Sebastian Weichwald. Beware of the simulated dag! causal discovery
598 benchmarks may be easy to game. *Advances in Neural Information Processing Systems*, 34:27772–27784,
599 2021.
- 600
- 601 Karen Sachs, Omar Perez, Dana Pe’er, Douglas A Lauffenburger, and Garry P Nolan. Causal protein-signaling
602 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
- 597 Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-Gaussian acyclic model
598 for causal discovery. *Journal of Machine Learning Research*, 7(Oct):2003–2030, 2006.
- 599
- 600 Tomi Silander and Petri Myllymäki. A simple approach for finding the globally optimal Bayesian network
601 structure. In *Conference on Uncertainty in Artificial Intelligence*, 2006.
- 602
- 603 Ajit P. Singh and Andrew W. Moore. Finding optimal Bayesian networks by dynamic programming. Technical
604 report, Carnegie Mellon University, 2005.
- 605
- 606 Daniel A. Spielman. Lecture notes in spectral graph theory: The adjacency matrix and the nth eigenvalue, September
607 2012. URL <https://www.cs.yale.edu/homes/spielman/561/2012/lect03-12.pdf>.
- 608
- 609 Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social science
610 computer review*, 9(1):62–72, 1991.
- 611
- 612 Peter Spirtes, Chris Meek, and Thomas Richardson. Causal inference in the presence of latent variables and
613 selection bias. In *Conference on Uncertainty in Artificial Intelligence*, 1995.
- 614
- 615 Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*.
616 MIT press, 2000.
- 617
- 618 Matthew J Vowels, Necati Cihan Camgoz, and Richard Bowden. D’ya like dags? a survey on structure learning
619 and causal discovery. *arXiv preprint arXiv:2103.02582*, 2021.
- 620
- 621 Dennis Wei, Tian Gao, and Yue Yu. Dags with No Fears: A closer look at continuous optimization for learning
622 bayesian networks. *Advances in Neural Information Processing Systems*, 33:3895–3906, 2020.
- 623
- 624 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. DAG-GNN: Dag structure learning with graph neural networks. In
625 *International Conference on Machine Learning*, pages 7154–7163. PMLR, 2019.
- 626
- 627 Yue Yu, Tian Gao, Naiyu Yin, and Qiang Ji. DAGs with no curl: An efficient dag structure learning approach. In
628 *Proceedings of the 38th International Conference on Machine Learning*, 2021.
- 629
- 630 Changhe Yuan and Brandon Malone. Learning optimal Bayesian networks: A shortest path perspective. *Journal
of Artificial Intelligence Research*, 48(1):23–65, 2013.
- 631
- 632 Changhe Yuan, Brandon Malone, and Xiaojian Wu. Learning optimal Bayesian networks using A* search. In
633 *International Joint Conference on Artificial Intelligence*, 2011.
- 634
- 635 Bin Zhang, Chris Gaiteri, Liviu-Gabriel Bodea, Zhi Wang, Joshua McElwee, Alexei A Podtelezhnikov, Chun-
636 sheng Zhang, Tao Xie, Linh Tran, Radu Dobrin, et al. Integrated systems approach identifies genetic nodes
637 and networks in late-onset alzheimer’s disease. *Cell*, 153(3):707–720, 2013.
- 638
- 639 Zhen Zhang, Ignavier Ng, Dong Gong, Yuhang Liu, Ehsan Abbasnejad, Mingming Gong, Kun Zhang, and
640 Javen Qinfeng Shi. Truncated matrix power iteration for differentiable dag learning. *Advances in Neural
Information Processing Systems*, 35:18390–18402, 2022.
- 641
- 642 Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. DAGs with No Tears: Continuous
643 optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.
- 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). \square

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). \square

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

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

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

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

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 $\{\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d} | \rho(\tilde{\mathbf{B}}) < r\}$, 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

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

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

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

770 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
 771 proved the proposition. \square

773 A.5 PROOF OF PROPOSITION 4

774 **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-
 775 jacency matrix $\tilde{\mathbf{B}} \in \mathbb{R}_{\geq 0}^{d \times d}$ with spectral radius $\rho(\tilde{\mathbf{B}}) \leq \min(\lim_{i \rightarrow \infty} c_i^1/c_{i+1}^1, \lim_{i \rightarrow \infty} c_i^2/c_{i+1}^2)$,
 776 the following three statements are equivalent:

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

780 *Proof.* By properties of analytic functions, we have

$$781 \quad f_1(x) + f_2(x) = c_0^1 + c_0^2 + \sum_{i=1}^{\infty} (c_i^1 + c_i^2)x^i \tag{30}$$

782 is an analytic function, and its convergence radius is given by

$$783 \quad \lim_{i \rightarrow \infty} (c_i^1 + c_i^2)/(c_{i+1}^1 + c_{i+1}^2) = \min(\lim_{i \rightarrow \infty} c_i^1/c_{i+1}^1, \lim_{i \rightarrow \infty} c_i^2/c_{i+1}^2), \tag{31}$$

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

786 A.6 PROOF OF PROPOSITION 5

787 **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
 788 and only if

$$789 \quad \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n} = d,$$

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

$$791 \quad \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n}\| \leq \|\nabla_{\tilde{\mathbf{B}}} \text{tr}(\mathbf{I} - \tilde{\mathbf{B}}/s)^{-n-k}\|,$$

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

794 *Proof.* By Proposition 4 or Proposition 3, it would be straightforward that $\text{tr}(\mathbf{I} - \tilde{\mathbf{B}})^{-n} = d$ is a
 795 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
 796 DAG.

For the norm of gradients, it is straightforward that

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

For arbitrary n we have

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

and obviously the coefficients are 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 $j > 0, k > 0$

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

As a result, we have

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

□

A.7 PROOF OF PROPOSITION 6

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

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

where \otimes denotes the Kronecker product.

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

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

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

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

□

A.8 PROOF OF PROPOSITION 7

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*

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

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

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

$$872 \quad \geq \max_{\mathbf{u}: \mathbf{u} \geq 0, \|\mathbf{u}\|_2=1} \mathbf{u}^\top \nabla_{\tilde{\mathbf{B}}}^2 \text{tr}f_2(\tilde{\mathbf{B}}) \mathbf{u} = \rho(\nabla_{\tilde{\mathbf{B}}}^2 \text{tr}f_2(\tilde{\mathbf{B}}))$$

874 \square
 875

876 B HYPER PARAMETERS 877

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.

885 C EXTRA EXPERIMENTAL RESULTS 886

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 **Table 5:** DAG learning performance (measured in true positive rate, the higher the better, best
930 results in **bold**) of different algorithms on large scale (500-2000 nodes) graphs with different noise
931 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	0.95 ± 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 ± 0.01	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

962
963
964
965
966
967
968
969
970
971

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983

984 Table 6: DAG learning performance (measured in false detection rate, the lower the better, best
 985 results in **bold**) of different algorithms on large scale (500-2000 nodes) graphs with different noise
 986 distributions. Our algorithm performs better than previous approaches.
 987

	Graphs	Nodes	DAGMA	Order-1	Order-2	Order-3	Order-4
988	ER2-gauss	500	0.02 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
989	ER2-gauss	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
990	ER2-gauss	2000	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.01	0.02 ± 0.00
991	ER3-gauss	500	0.04 ± 0.02	0.04 ± 0.02	0.03 ± 0.02	0.03 ± 0.01	0.04 ± 0.02
992	ER3-gauss	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
993	ER3-gauss	2000	0.06 ± 0.01	0.06 ± 0.02	0.05 ± 0.01	0.04 ± 0.01	0.04 ± 0.01
994	ER4-gauss	500	0.08 ± 0.03	0.08 ± 0.03	0.08 ± 0.04	0.07 ± 0.04	0.07 ± 0.03
995	ER4-gauss	1000	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.02	0.09 ± 0.03	0.10 ± 0.03
996	ER4-gauss	2000	0.14 ± 0.01	0.13 ± 0.01	0.13 ± 0.02	0.12 ± 0.02	0.12 ± 0.01
997	ER2-exp	500	0.03 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.02	0.02 ± 0.02
998	ER2-exp	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
999	ER3-exp	500	0.05 ± 0.02	0.04 ± 0.02	0.03 ± 0.01	0.04 ± 0.02	0.04 ± 0.02
1000	ER3-exp	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
1001	ER4-exp	500	0.09 ± 0.04	0.09 ± 0.04	0.09 ± 0.05	0.07 ± 0.04	0.06 ± 0.03
1002	ER4-exp	1000	0.12 ± 0.03	0.11 ± 0.02	0.11 ± 0.02	0.09 ± 0.03	0.10 ± 0.02
1003	ER2-gumbel	500	0.03 ± 0.02	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
1004	ER2-gumbel	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.00	0.01 ± 0.01	0.02 ± 0.01
1005	ER2-gumbel	2000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.00	0.01 ± 0.01	0.02 ± 0.01
1006	ER3-gumbel	500	0.06 ± 0.02	0.04 ± 0.02	0.03 ± 0.03	0.03 ± 0.01	0.04 ± 0.02
1007	ER3-gumbel	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.02	0.04 ± 0.01	0.03 ± 0.01
1008	ER4-gumbel	500	0.10 ± 0.04	0.09 ± 0.05	0.09 ± 0.04	0.08 ± 0.04	0.08 ± 0.05
1009	ER4-gumbel	1000	0.14 ± 0.03	0.12 ± 0.03	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.03
1010	SF2-gauss	500	0.02 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
1011	SF2-gauss	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
1012	SF2-gauss	2000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
1013	SF3-gauss	500	0.04 ± 0.02	0.04 ± 0.02	0.03 ± 0.02	0.03 ± 0.01	0.04 ± 0.02
1014	SF3-gauss	1000	0.06 ± 0.01	0.06 ± 0.02	0.05 ± 0.01	0.04 ± 0.01	0.04 ± 0.01
1015	SF4-gauss	500	0.08 ± 0.03	0.08 ± 0.03	0.08 ± 0.04	0.07 ± 0.04	0.07 ± 0.03
1016	SF4-gauss	1000	0.12 ± 0.03	0.11 ± 0.03	0.10 ± 0.02	0.09 ± 0.03	0.10 ± 0.03
1017	SF4-gauss	2000	0.14 ± 0.01	0.13 ± 0.01	0.13 ± 0.02	0.12 ± 0.02	0.12 ± 0.01
1018	SF2-exp	500	0.03 ± 0.02	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.02	0.02 ± 0.02
1019	SF2-exp	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
1020	SF3-exp	500	0.05 ± 0.02	0.04 ± 0.02	0.03 ± 0.01	0.04 ± 0.02	0.04 ± 0.02
1021	SF3-exp	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.01	0.03 ± 0.01	0.03 ± 0.01
1022	SF4-exp	500	0.09 ± 0.04	0.09 ± 0.04	0.09 ± 0.05	0.07 ± 0.04	0.06 ± 0.03
1023	SF4-exp	1000	0.12 ± 0.03	0.11 ± 0.02	0.11 ± 0.02	0.09 ± 0.03	0.10 ± 0.02
1024	SF2-gumbel	500	0.03 ± 0.02	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01	0.01 ± 0.01
1025	SF2-gumbel	1000	0.02 ± 0.01	0.02 ± 0.01	0.01 ± 0.00	0.01 ± 0.01	0.02 ± 0.01
1026	SF3-gumbel	500	0.06 ± 0.02	0.04 ± 0.02	0.03 ± 0.03	0.03 ± 0.01	0.04 ± 0.02
1027	SF3-gumbel	1000	0.06 ± 0.01	0.05 ± 0.01	0.04 ± 0.02	0.04 ± 0.01	0.03 ± 0.01
1028	SF4-gumbel	500	0.10 ± 0.04	0.09 ± 0.05	0.09 ± 0.04	0.08 ± 0.04	0.08 ± 0.05
1029	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 Table 7: DAG learning performance (measured in running time (seconds), the lower the better, best
1039 results in **bold**) of different algorithms on large scale (500-2000 nodes) graphs with different noise
1040 distributions. Our algorithm performs better than previous approaches.

1041

	Graphs	Nodes	DAGMA	Order-1	Order-2	Order-3	Order-4
1042	ER2-gauss	500	171.79 ± 18.30	294.19 ± 62.47	493.73 ± 28.79	482.01 ± 17.59	512.86 ± 33.87
1043	ER2-gauss	1000	364.95 ± 38.47	687.38 ± 69.45	1261.33 ± 47.66	1256.89 ± 33.76	1329.92 ± 315.91
1044	ER2-gauss	2000	1187.17 ± 108.06	3515.56 ± 310.54	6063.65 ± 554.57	6300.09 ± 197.69	6360.65 ± 321.65
1045	ER3-gauss	500	238.54 ± 58.34	394.95 ± 110.87	516.49 ± 14.45	524.93 ± 48.79	528.08 ± 49.59
1046	ER3-gauss	1000	501.24 ± 59.30	1004.85 ± 162.25	1365.24 ± 21.10	1355.23 ± 139.98	1349.25 ± 84.01
1047	ER3-gauss	2000	1636.23 ± 101.64	4903.56 ± 782.06	6037.67 ± 765.47	7072.04 ± 2168.01	6969.54 ± 724.52
1048	ER4-gauss	500	347.06 ± 55.40	519.77 ± 56.88	532.34 ± 10.16	517.46 ± 12.95	534.98 ± 52.22
1049	ER4-gauss	1000	798.01 ± 27.84	1328.32 ± 48.32	1318.45 ± 62.71	1348.08 ± 15.79	1312.53 ± 138.14
1050	ER4-gauss	2000	2461.67 ± 1.49	6520.96 ± 311.23	6560.10 ± 13.55	7666.47 ± 2151.48	7067.87 ± 1011.11
1051	ER2-exp	500	167.08 ± 25.75	291.87 ± 63.66	492.14 ± 30.49	496.83 ± 20.55	504.17 ± 22.59
1052	ER2-exp	1000	359.97 ± 28.62	708.23 ± 61.00	1245.71 ± 71.43	1279.88 ± 41.63	1277.89 ± 45.39
1053	ER3-exp	500	235.71 ± 56.44	440.70 ± 189.91	564.88 ± 145.09	559.54 ± 144.42	550.44 ± 129.03
1054	ER3-exp	1000	515.88 ± 83.18	1100.18 ± 226.98	1371.65 ± 136.38	1401.29 ± 283.39	1503.85 ± 321.35
1055	ER4-exp	500	358.94 ± 46.65	510.53 ± 46.91	513.70 ± 49.38	522.16 ± 15.09	508.01 ± 34.87
1056	ER4-exp	1000	778.40 ± 51.69	1344.07 ± 26.71	1324.00 ± 111.40	1347.33 ± 21.06	1298.91 ± 100.30
1057	ER2-gumbel	500	161.36 ± 24.90	255.86 ± 33.56	501.55 ± 85.14	490.99 ± 11.06	501.44 ± 23.26
1058	ER2-gumbel	1000	330.53 ± 39.10	656.98 ± 62.61	1245.11 ± 47.50	1276.58 ± 15.41	1266.35 ± 38.53
1059	ER3-gumbel	500	232.45 ± 50.71	381.03 ± 85.67	521.98 ± 8.54	514.79 ± 12.92	506.95 ± 18.01
1060	ER3-gumbel	1000	525.92 ± 93.20	1013.67 ± 168.84	1331.88 ± 99.25	1302.93 ± 101.85	1369.44 ± 39.47
1061	ER4-gumbel	500	366.06 ± 33.40	514.95 ± 57.57	540.93 ± 17.35	530.18 ± 20.49	519.87 ± 13.98
1062	ER4-gumbel	1000	805.91 ± 31.96	1367.02 ± 36.36	1260.93 ± 137.27	1434.41 ± 154.30	1335.20 ± 61.78
1063	SF2-gauss	500	171.79 ± 18.30	294.19 ± 62.47	493.73 ± 28.79	482.01 ± 17.59	512.86 ± 33.87
1064	SF2-gauss	1000	364.95 ± 38.47	687.38 ± 69.45	1261.33 ± 47.66	1256.89 ± 33.76	1329.92 ± 315.91
1065	SF2-gauss	2000	1187.17 ± 108.06	3515.56 ± 310.54	6063.65 ± 554.57	6300.09 ± 197.69	6360.65 ± 321.65
1066	SF3-gauss	500	238.54 ± 58.34	394.95 ± 110.87	516.49 ± 14.45	524.93 ± 48.79	528.08 ± 49.59
1067	SF3-gauss	1000	501.24 ± 59.30	1004.85 ± 162.25	1365.24 ± 21.10	1355.23 ± 139.98	1349.25 ± 84.01
1068	SF3-gauss	2000	1636.23 ± 101.64	4903.56 ± 782.06	6037.67 ± 765.47	7072.04 ± 2168.01	6969.54 ± 724.52
1069	SF4-gauss	500	347.06 ± 55.40	519.77 ± 56.88	532.34 ± 10.16	517.46 ± 12.95	534.98 ± 52.22
1070	SF4-gauss	1000	798.01 ± 27.84	1328.32 ± 48.32	1318.45 ± 62.71	1348.08 ± 15.79	1312.53 ± 138.14
1071	SF4-gauss	2000	2461.67 ± 1.49	6520.96 ± 311.23	6560.10 ± 13.55	7666.47 ± 2151.48	7067.87 ± 1011.11
1072	SF2-exp	500	167.08 ± 25.75	291.87 ± 63.66	492.14 ± 30.49	496.83 ± 20.55	504.17 ± 22.59
1073	SF2-exp	1000	359.97 ± 28.62	708.23 ± 61.00	1245.71 ± 71.43	1279.88 ± 41.63	1277.89 ± 45.39
1074	SF3-exp	500	235.71 ± 56.44	440.70 ± 189.91	564.88 ± 145.09	559.54 ± 144.42	550.44 ± 129.03
1075	SF3-exp	1000	515.88 ± 83.18	1100.18 ± 226.98	1371.65 ± 136.38	1401.29 ± 283.39	1503.85 ± 321.35
1076	SF4-exp	500	358.94 ± 46.65	510.53 ± 46.91	513.70 ± 49.38	522.16 ± 15.09	508.01 ± 34.87
1077	SF4-exp	1000	778.40 ± 51.69	1344.07 ± 26.71	1324.00 ± 111.40	1347.33 ± 21.06	1298.91 ± 100.30
1078	SF2-gumbel	500	161.36 ± 24.90	255.86 ± 33.56	501.55 ± 85.14	490.99 ± 11.06	501.44 ± 23.26
1079	SF2-gumbel	1000	330.53 ± 39.10	656.98 ± 62.61	1245.11 ± 47.50	1276.58 ± 15.41	1266.35 ± 38.53
1080	SF3-gumbel	500	232.45 ± 50.71	381.03 ± 85.67	521.98 ± 8.54	514.79 ± 12.92	506.95 ± 18.01
1081	SF3-gumbel	1000	525.92 ± 93.20	1013.67 ± 168.84	1331.88 ± 99.25	1302.93 ± 101.85	1369.44 ± 39.47
1082	SF4-gumbel	500	366.06 ± 33.40	514.95 ± 57.57	540.93 ± 17.35	530.18 ± 20.49	519.87 ± 13.98
1083	SF4-gumbel	1000	805.91 ± 31.96	1367.02 ± 36.36	1260.93 ± 137.27	1434.41 ± 154.30	1335.20 ± 61.78

Table 8: DAG learning performance (measured in structural hamming distance, the lower the better, best results in **bold**) of different algorithms on small scale (10-100 nodes) ER{2,3,4} graphs with different noise distributions. Our algorithm performs better than previous approaches.

	Graphs	Nodes	MMPC	GES	NOTEARS	DAGMA	Order-1	Order-2	Order-3	Order-4
1080	ER2-gauss	10	21.01 \pm 3.77	14.66 \pm 9.25	3.01 \pm 2.67	2.01 \pm 1.91	0.74 \pm 1.22	0.93 \pm 1.49	0.97 \pm 1.65	1.27 \pm 2.11
1081	ER2-gauss	30	65.23 \pm 9.26	57.95 \pm 44.65	6.09 \pm 6.31	3.76 \pm 3.70	1.88 \pm 2.97	1.57 \pm 2.29	1.49 \pm 2.32	1.87 \pm 2.89
1082	ER2-gauss	50	106.74 \pm 11.35	90.49 \pm 58.75	11.59 \pm 10.25	4.21 \pm 3.66	2.28 \pm 3.06	2.05 \pm 2.96	2.21 \pm 2.98	2.43 \pm 3.69
1083	ER2-gauss	100	215.36 \pm 14.83	155.61 \pm 76.87	22.58 \pm 17.69	7.45 \pm 6.70	3.49 \pm 3.97	3.64 \pm 4.48	4.09 \pm 4.40	4.01 \pm 4.54
1084	ER3-gauss	10	32.04 \pm 3.39	26.39 \pm 6.85	8.91 \pm 4.33	6.96 \pm 3.69	3.36 \pm 2.91	3.35 \pm 2.97	3.21 \pm 2.83	3.51 \pm 2.83
1085	ER3-gauss	30	97.56 \pm 9.71	158.57 \pm 49.92	17.69 \pm 12.29	9.22 \pm 6.24	4.83 \pm 4.56	5.33 \pm 5.52	5.54 \pm 5.71	4.96 \pm 5.94
1086	ER3-gauss	50	160.21 \pm 12.67	285.88 \pm 96.26	30.70 \pm 18.48	13.05 \pm 8.68	6.60 \pm 6.56	6.43 \pm 5.72	6.55 \pm 6.09	6.44 \pm 5.54
1087	ER3-gauss	100	321.30 \pm 18.01	506.91 \pm 202.63	64.27 \pm 33.28	21.66 \pm 15.48	11.59 \pm 9.90	10.69 \pm 9.92	10.86 \pm 9.02	11.28 \pm 9.99
1088	ER4-gauss	10	40.37 \pm 2.24	29.07 \pm 5.77	13.71 \pm 3.93	11.57 \pm 3.27	6.89 \pm 3.39	6.75 \pm 3.15	7.08 \pm 3.40	7.04 \pm 3.31
1089	ER4-gauss	30	126.27 \pm 10.70	222.49 \pm 38.99	39.18 \pm 21.79	20.67 \pm 9.26	11.12 \pm 7.76	10.54 \pm 7.35	12.58 \pm 8.81	12.81 \pm 10.74
1090	ER4-gauss	50	210.43 \pm 13.13	495.14 \pm 81.60	60.61 \pm 25.25	24.96 \pm 12.48	14.39 \pm 9.72	12.89 \pm 9.52	15.15 \pm 10.31	16.22 \pm 10.17
1091	ER4-gauss	100	424.49 \pm 20.13	1047.05 \pm 250.05	118.16 \pm 51.53	42.81 \pm 23.74	28.31 \pm 25.58	24.47 \pm 24.31	23.22 \pm 23.60	23.03 \pm 21.51
1092	ER2-exp	10	20.88 \pm 3.70	15.23 \pm 9.49	3.05 \pm 2.64	2.19 \pm 1.87	0.80 \pm 1.34	0.90 \pm 1.57	1.03 \pm 1.74	1.35 \pm 2.30
1093	ER2-exp	30	64.98 \pm 9.21	59.37 \pm 46.71	6.59 \pm 7.08	3.83 \pm 4.10	2.22 \pm 3.63	1.85 \pm 2.91	1.40 \pm 2.37	2.09 \pm 3.28
1094	ER2-exp	50	106.77 \pm 10.85	95.43 \pm 56.48	11.51 \pm 9.70	4.21 \pm 3.43	2.12 \pm 2.75	1.97 \pm 2.47	2.62 \pm 3.01	2.80 \pm 3.53
1095	ER2-exp	100	215.70 \pm 14.80	159.55 \pm 77.19	21.87 \pm 17.48	7.31 \pm 6.17	3.71 \pm 4.77	4.11 \pm 4.58	4.54 \pm 4.72	4.82 \pm 5.15
1096	ER3-exp	10	31.90 \pm 3.41	26.29 \pm 6.80	9.32 \pm 4.42	6.98 \pm 3.43	3.57 \pm 2.91	3.43 \pm 2.96	3.10 \pm 2.89	3.68 \pm 2.82
1097	ER3-exp	30	97.74 \pm 9.81	154.20 \pm 48.45	17.54 \pm 11.51	9.11 \pm 5.51	5.07 \pm 4.47	5.51 \pm 5.51	5.60 \pm 5.73	4.81 \pm 5.76
1098	ER3-exp	50	160.16 \pm 12.49	288.43 \pm 99.00	31.32 \pm 19.96	14.16 \pm 9.73	7.92 \pm 7.38	7.63 \pm 5.94	7.15 \pm 6.32	7.18 \pm 5.69
1099	ER3-exp	100	321.66 \pm 17.89	494.33 \pm 188.84	66.39 \pm 30.84	22.02 \pm 15.03	12.26 \pm 9.06	12.44 \pm 11.03	11.52 \pm 9.37	11.14 \pm 8.36
1100	ER4-exp	10	40.44 \pm 2.15	28.79 \pm 5.56	13.92 \pm 3.89	11.85 \pm 3.42	7.01 \pm 3.44	6.93 \pm 3.36	7.07 \pm 3.55	7.10 \pm 3.32
1101	ER4-exp	30	125.98 \pm 10.62	221.88 \pm 39.40	36.55 \pm 18.74	21.00 \pm 10.26	10.73 \pm 7.83	11.55 \pm 7.70	13.19 \pm 9.95	12.87 \pm 11.09
1102	ER4-exp	50	210.42 \pm 13.43	494.91 \pm 77.20	60.90 \pm 25.93	27.10 \pm 12.94	14.76 \pm 9.50	13.83 \pm 9.46	14.92 \pm 10.98	16.44 \pm 12.95
1103	ER4-exp	100	424.65 \pm 19.82	1043.47 \pm 239.69	116.91 \pm 48.63	42.64 \pm 23.78	27.02 \pm 24.07	24.32 \pm 22.06	23.39 \pm 22.89	24.12 \pm 22.40
1104	ER2-gumbel	10	20.89 \pm 3.72	16.07 \pm 8.92	1.69 \pm 2.04	1.11 \pm 1.52	0.67 \pm 1.36	0.82 \pm 1.60	0.79 \pm 1.65	1.00 \pm 2.24
1105	ER2-gumbel	30	64.88 \pm 9.55	61.33 \pm 43.70	5.65 \pm 7.35	2.10 \pm 2.92	1.76 \pm 2.79	1.69 \pm 2.90	2.13 \pm 3.06	2.28 \pm 3.74
1106	ER2-gumbel	50	106.77 \pm 10.76	92.35 \pm 57.81	10.24 \pm 10.44	2.81 \pm 3.30	2.43 \pm 3.03	2.64 \pm 3.81	3.41 \pm 4.63	3.77 \pm 4.56
1107	ER2-gumbel	100	215.52 \pm 15.57	161.75 \pm 89.91	22.39 \pm 19.17	4.89 \pm 5.94	4.21 \pm 5.17	4.78 \pm 5.64	5.59 \pm 5.42	6.63 \pm 6.00
1108	ER3-gumbel	10	31.95 \pm 3.37	26.29 \pm 6.69	5.93 \pm 3.95	4.05 \pm 2.96	1.90 \pm 2.39	2.44 \pm 2.71	2.09 \pm 2.31	1.75 \pm 2.23
1109	ER3-gumbel	30	97.72 \pm 9.72	158.60 \pm 49.62	13.44 \pm 11.65	4.97 \pm 4.07	3.72 \pm 4.11	4.13 \pm 4.29	5.99 \pm 6.27	5.96 \pm 6.92
1110	ER3-gumbel	50	160.35 \pm 12.66	281.98 \pm 102.11	27.28 \pm 20.25	9.43 \pm 7.56	7.59 \pm 5.76	7.89 \pm 6.40	8.17 \pm 8.16	
1111	ER3-gumbel	100	321.87 \pm 18.02	496.90 \pm 193.46	60.12 \pm 29.08	15.06 \pm 13.15	12.62 \pm 12.83	12.46 \pm 12.84	13.87 \pm 10.60	
1112	ER4-gumbel	10	40.48 \pm 2.13	28.69 \pm 5.93	10.24 \pm 4.08	7.75 \pm 2.94	3.96 \pm 3.11	4.55 \pm 3.54	4.92 \pm 3.38	4.90 \pm 3.46
1113	ER4-gumbel	30	126.04 \pm 10.54	219.39 \pm 43.15	29.45 \pm 19.08	12.35 \pm 6.76	8.14 \pm 7.53	11.29 \pm 9.49	11.84 \pm 9.91	
1114	ER4-gumbel	50	210.16 \pm 13.34	499.02 \pm 78.43	55.57 \pm 28.67	18.03 \pm 15.20	12.07 \pm 9.36	10.92 \pm 9.04	13.44 \pm 10.94	15.06 \pm 11.37
1115	ER4-gumbel	100	424.19 \pm 20.06	1031.32 \pm 243.38	114.40 \pm 52.86	29.94 \pm 22.69	27.00 \pm 28.51	22.36 \pm 24.07	21.24 \pm 23.58	25.37 \pm 25.67
1116	SF2-gauss	10	14.82 \pm 1.99	5.01 \pm 6.58	3.01 \pm 2.67	2.01 \pm 1.91	0.74 \pm 1.22	0.93 \pm 1.49	0.97 \pm 1.65	1.27 \pm 2.11
1117	SF2-gauss	30	54.01 \pm 3.54	19.59 \pm 19.48	6.09 \pm 6.31	3.76 \pm 3.70	1.88 \pm 2.97	1.49 \pm 2.32	1.87 \pm 2.89	
1118	SF2-gauss	50	97.30 \pm 5.13	43.39 \pm 40.68	11.59 \pm 10.25	4.21 \pm 3.66	2.28 \pm 3.06	2.05 \pm 2.96	2.21 \pm 2.98	2.43 \pm 3.69
1119	SF2-gauss	100	215.89 \pm 8.24	106.51 \pm 62.29	22.58 \pm 17.69	7.45 \pm 6.70	3.49 \pm 3.97	3.64 \pm 4.48	4.09 \pm 4.40	4.01 \pm 4.54
1120	SF3-gauss	10	16.61 \pm 2.77	7.78 \pm 8.34	8.91 \pm 4.33	6.96 \pm 3.69	3.36 \pm 2.91	3.35 \pm 2.97	3.21 \pm 2.83	3.51 \pm 2.83
1121	SF3-gauss	30	65.90 \pm 7.29	31.51 \pm 32.61	17.69 \pm 12.29	9.22 \pm 6.24	4.83 \pm 4.56	5.33 \pm 5.52	5.54 \pm 5.71	4.96 \pm 5.94
1122	SF3-gauss	50	119.92 \pm 10.81	69.43 \pm 58.72	30.70 \pm 18.48	13.05 \pm 8.68	6.60 \pm 6.56	6.43 \pm 5.72	6.55 \pm 6.09	6.44 \pm 5.54
1123	SF3-gauss	100	271.55 \pm 16.84	157.11 \pm 99.22	64.27 \pm 33.28	21.66 \pm 15.48	11.59 \pm 9.90	10.69 \pm 9.92	10.86 \pm 9.02	11.28 \pm 9.99
1124	SF4-gauss	10	17.53 \pm 3.20	6.54 \pm 6.84	13.71 \pm 3.93	11.57 \pm 3.27	6.89 \pm 3.39	6.75 \pm 3.15	7.08 \pm 3.40	7.04 \pm 3.31
1125	SF4-gauss	30	73.81 \pm 9.41	46.72 \pm 41.14	39.18 \pm 21.79	20.67 \pm 9.26	11.12 \pm 7.76	10.54 \pm 7.35	12.58 \pm 8.81	12.81 \pm 10.74
1126	SF4-gauss	50	137.96 \pm 14.04	84.43 \pm 63.15	60.99 \pm 25.93	24.96 \pm 12.48	14.39 \pm 9.72	12.89 \pm 9.52	15.15 \pm 10.31	16.22 \pm 10.17
1127	SF4-gauss	100	314.36 \pm 24.83	176.43 \pm 114.42	116.91 \pm 48.63	42.64 \pm 23.78	28.31 \pm 25.58	24.47 \pm 24.31	23.22 \pm 23.60	23.03 \pm 21.51
1128	SF2-exp	10	14.72 \pm 1.88	5.18 \pm 6.51	3.05 \pm 2.64	2.19 \pm 1.87	0.80 \pm 1.34	0.90 \pm 1.57	1.03 \pm 1.74	1.35 \pm 2.30
1129	SF2-exp	30	54.19 \pm 3.69	23.86 \pm 24.42	6.59 \pm 7.08	3.83 \pm 4.10	2.22 \pm 3.63	1.85 \pm 2.91	1.40 \pm 2.37	2.09 \pm 3.28
1130	SF2-exp	50	97.15 \pm 5.42	47.29 \pm 42.59	11.51 \pm 9.70	4.21 \pm 3.43	2.12 \pm 2.75	1.97 \pm 2.47	2.62 \pm 3.01	2.80 \pm 3.53
1131	SF2-exp	100	216.20 \pm 7.98	113.30 \pm 63.41	21.87 \pm 17.48	7.31 \pm 6.17	3.71 \pm 4.77	4.11 \pm 4.58	4.54 \pm 4.72	4.82 \pm 5.15
1132	SF3-exp	10	16.58 \pm 2.77	6.54 \pm 7.60	9.32 $\pm</math$					

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

Table 9: DAG learning performance (measured in true positive rate, the higher the better, best results in **bold**) of different algorithms on small scale (10-100 nodes) ER{2,3,4} graphs with different noise distributions. Our algorithm performs better than previous approaches.

Graphs	Nodes	MMPC	GES	NOTEARS	DAGMA	Order-1	Order-2	Order-3	Order-4
ER2-gauss	10	0.59 ± 0.14	0.70 ± 0.23	0.88 ± 0.09	0.91 ± 0.07	0.97 ± 0.04	0.97 ± 0.04	0.97 ± 0.04	0.96 ± 0.05
ER2-gauss	30	0.53 ± 0.11	0.80 ± 0.14	0.93 ± 0.05	0.95 ± 0.04	0.98 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.98 ± 0.02
ER2-gauss	50	0.55 ± 0.09	0.83 ± 0.10	0.93 ± 0.05	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER2-gauss	100	0.58 ± 0.07	0.87 ± 0.06	0.94 ± 0.04	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER3-gauss	10	0.34 ± 0.07	0.48 ± 0.15	0.75 ± 0.11	0.80 ± 0.10	0.91 ± 0.08	0.91 ± 0.08	0.91 ± 0.07	0.91 ± 0.07
ER3-gauss	30	0.29 ± 0.07	0.63 ± 0.13	0.88 ± 0.06	0.92 ± 0.04	0.97 ± 0.02	0.97 ± 0.03	0.97 ± 0.03	0.97 ± 0.03
ER3-gauss	50	0.30 ± 0.06	0.69 ± 0.10	0.89 ± 0.06	0.93 ± 0.03	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.02
ER3-gauss	100	0.33 ± 0.04	0.78 ± 0.07	0.89 ± 0.04	0.95 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
ER4-gauss	10	0.28 ± 0.05	0.43 ± 0.14	0.68 ± 0.09	0.73 ± 0.07	0.84 ± 0.07	0.85 ± 0.07	0.84 ± 0.07	0.84 ± 0.07
ER4-gauss	30	0.18 ± 0.04	0.54 ± 0.10	0.80 ± 0.08	0.87 ± 0.04	0.94 ± 0.03	0.95 ± 0.03	0.94 ± 0.03	0.94 ± 0.03
ER4-gauss	50	0.18 ± 0.03	0.57 ± 0.09	0.83 ± 0.05	0.91 ± 0.03	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02
ER4-gauss	100	0.20 ± 0.03	0.68 ± 0.07	0.85 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
ER2-exp	10	0.58 ± 0.14	0.67 ± 0.24	0.87 ± 0.09	0.91 ± 0.07	0.97 ± 0.04	0.97 ± 0.04	0.97 ± 0.04	0.96 ± 0.05
ER2-exp	30	0.54 ± 0.11	0.80 ± 0.14	0.93 ± 0.05	0.95 ± 0.04	0.98 ± 0.02	0.98 ± 0.02	0.99 ± 0.02	0.98 ± 0.02
ER2-exp	50	0.56 ± 0.10	0.82 ± 0.09	0.93 ± 0.05	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER2-exp	100	0.59 ± 0.07	0.87 ± 0.06	0.94 ± 0.04	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER3-exp	10	0.34 ± 0.08	0.48 ± 0.16	0.75 ± 0.11	0.81 ± 0.09	0.90 ± 0.07	0.91 ± 0.08	0.92 ± 0.07	0.90 ± 0.07
ER3-exp	30	0.29 ± 0.07	0.64 ± 0.13	0.88 ± 0.06	0.92 ± 0.04	0.97 ± 0.02	0.96 ± 0.03	0.97 ± 0.03	0.97 ± 0.03
ER3-exp	50	0.31 ± 0.06	0.69 ± 0.10	0.88 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.98 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
ER3-exp	100	0.33 ± 0.04	0.78 ± 0.07	0.89 ± 0.04	0.95 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
ER4-exp	10	0.28 ± 0.05	0.44 ± 0.13	0.68 ± 0.09	0.72 ± 0.07	0.84 ± 0.08	0.84 ± 0.07	0.84 ± 0.08	0.84 ± 0.07
ER4-exp	30	0.17 ± 0.04	0.54 ± 0.11	0.81 ± 0.07	0.87 ± 0.05	0.94 ± 0.03	0.94 ± 0.03	0.94 ± 0.03	0.94 ± 0.04
ER4-exp	50	0.18 ± 0.03	0.58 ± 0.09	0.83 ± 0.05	0.90 ± 0.03	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02
ER4-exp	100	0.20 ± 0.03	0.68 ± 0.07	0.85 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
ER2-gumbel	10	0.57 ± 0.15	0.64 ± 0.24	0.93 ± 0.08	0.96 ± 0.05	0.98 ± 0.04	0.98 ± 0.04	0.98 ± 0.04	0.98 ± 0.04
ER2-gumbel	30	0.54 ± 0.11	0.78 ± 0.13	0.95 ± 0.05	0.98 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.99 ± 0.02
ER2-gumbel	50	0.56 ± 0.10	0.83 ± 0.09	0.95 ± 0.04	0.98 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER2-gumbel	100	0.58 ± 0.07	0.87 ± 0.06	0.95 ± 0.03	0.99 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER3-gumbel	10	0.35 ± 0.08	0.49 ± 0.16	0.84 ± 0.10	0.89 ± 0.07	0.95 ± 0.06	0.94 ± 0.06	0.95 ± 0.06	0.95 ± 0.06
ER3-gumbel	30	0.28 ± 0.07	0.63 ± 0.13	0.92 ± 0.06	0.96 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.97 ± 0.02	0.97 ± 0.02
ER3-gumbel	50	0.30 ± 0.06	0.69 ± 0.10	0.92 ± 0.05	0.96 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.02
ER3-gumbel	100	0.33 ± 0.04	0.78 ± 0.08	0.92 ± 0.03	0.97 ± 0.02	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
ER4-gumbel	10	0.28 ± 0.05	0.44 ± 0.13	0.77 ± 0.09	0.82 ± 0.07	0.91 ± 0.07	0.90 ± 0.07	0.89 ± 0.07	0.89 ± 0.08
ER4-gumbel	30	0.17 ± 0.04	0.54 ± 0.11	0.87 ± 0.07	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.96 ± 0.03	0.96 ± 0.03
ER4-gumbel	50	0.18 ± 0.03	0.57 ± 0.09	0.87 ± 0.05	0.95 ± 0.03	0.98 ± 0.02	0.98 ± 0.01	0.97 ± 0.02	0.97 ± 0.02
ER4-gumbel	100	0.20 ± 0.03	0.69 ± 0.07	0.88 ± 0.06	0.96 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.01	0.98 ± 0.01
SF2-gauss	10	0.78 ± 0.11	0.91 ± 0.19	0.88 ± 0.09	0.91 ± 0.07	0.97 ± 0.04	0.97 ± 0.04	0.97 ± 0.04	0.96 ± 0.05
SF2-gauss	30	0.62 ± 0.09	0.95 ± 0.08	0.93 ± 0.05	0.95 ± 0.04	0.98 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.98 ± 0.02
SF2-gauss	50	0.57 ± 0.08	0.92 ± 0.11	0.93 ± 0.05	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF2-gauss	100	0.54 ± 0.06	0.92 ± 0.08	0.94 ± 0.04	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF3-gauss	10	0.73 ± 0.13	0.86 ± 0.19	0.75 ± 0.11	0.80 ± 0.10	0.91 ± 0.08	0.91 ± 0.08	0.91 ± 0.07	0.91 ± 0.07
SF3-gauss	30	0.53 ± 0.10	0.92 ± 0.09	0.88 ± 0.06	0.92 ± 0.04	0.97 ± 0.02	0.97 ± 0.03	0.97 ± 0.03	0.97 ± 0.03
SF3-gauss	50	0.47 ± 0.07	0.90 ± 0.10	0.89 ± 0.06	0.93 ± 0.03	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.02
SF3-gauss	100	0.43 ± 0.05	0.90 ± 0.09	0.88 ± 0.04	0.95 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF4-gauss	10	0.71 ± 0.13	0.89 ± 0.15	0.68 ± 0.09	0.73 ± 0.07	0.84 ± 0.07	0.85 ± 0.07	0.84 ± 0.07	0.84 ± 0.07
SF4-gauss	30	0.49 ± 0.08	0.87 ± 0.12	0.80 ± 0.08	0.87 ± 0.04	0.94 ± 0.03	0.95 ± 0.03	0.94 ± 0.03	0.94 ± 0.03
SF4-gauss	50	0.42 ± 0.07	0.88 ± 0.10	0.83 ± 0.05	0.91 ± 0.03	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02
SF4-gauss	100	0.39 ± 0.05	0.90 ± 0.09	0.85 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
SF2-exp	10	0.79 ± 0.11	0.91 ± 0.18	0.87 ± 0.09	0.91 ± 0.07	0.97 ± 0.04	0.97 ± 0.04	0.97 ± 0.04	0.96 ± 0.05
SF2-exp	30	0.62 ± 0.09	0.93 ± 0.11	0.93 ± 0.05	0.95 ± 0.04	0.98 ± 0.02	0.98 ± 0.02	0.99 ± 0.02	0.98 ± 0.02
SF2-exp	50	0.57 ± 0.08	0.92 ± 0.11	0.93 ± 0.05	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF2-exp	100	0.54 ± 0.06	0.91 ± 0.08	0.94 ± 0.04	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF3-exp	10	0.74 ± 0.13	0.90 ± 0.17	0.75 ± 0.11	0.81 ± 0.09	0.90 ± 0.07	0.91 ± 0.08	0.92 ± 0.07	0.90 ± 0.07
SF3-exp	30	0.53 ± 0.10	0.92 ± 0.10	0.88 ± 0.06	0.92 ± 0.04	0.97 ± 0.02	0.96 ± 0.03	0.97 ± 0.03	0.97 ± 0.03
SF3-exp	50	0.47 ± 0.07	0.91 ± 0.09	0.88 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.98 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
SF3-exp	100	0.43 ± 0.05	0.89 ± 0.10	0.89 ± 0.04	0.95 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF4-exp	10	0.71 ± 0.13	0.88 ± 0.16	0.68 ± 0.09	0.72 ± 0.07	0.84 ± 0.08	0.84 ± 0.07	0.84 ± 0.08	0.84 ± 0.07
SF4-exp	30	0.49 ± 0.08	0.88 ± 0.12	0.81 ± 0.07	0.87 ± 0.05	0.94 ± 0.03	0.94 ± 0.03	0.94 ± 0.03	0.94 ± 0.04
SF4-exp	50	0.42 ± 0.07	0.89 ± 0.10	0.83 ± 0.05	0.90 ± 0.03	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02
SF4-exp	100	0.39 ± 0.05	0.90 ± 0.09	0.85 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
SF2-gumbel	10	0.78 ± 0.11	0.91 ± 0.18	0.93 ± 0.08	0.96 ± 0.05	0.98 ± 0.04	0.98 ± 0.04	0.98 ± 0.04	0.98 ± 0.04
SF2-gumbel	30	0.62 ± 0.09	0.94 ± 0.08	0.95 ± 0.05	0.98 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.99 ± 0.02
SF2-gumbel	50	0.57 ± 0.08	0.92 ± 0.10	0.95 ± 0.04	0.98 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF2-gumbel	100	0.53 ± 0.06	0.92 ± 0.08	0.95 ± 0.03	0.99 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF3-gumbel	10	0.73 ± 0.13	0.89 ± 0.17	0.84 ± 0.10	0.89 ± 0.07	0.95 ± 0.06	0.94 ± 0.06	0.94 ± 0.06	0.95 ± 0.06
SF3-gumbel	30	0.54 ± 0.10	0.92 ± 0.10	0.92 ± 0.06	0.96 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.97 ± 0.02	0.97 ± 0.02
SF3-gumbel	50	0.48 ± 0.07	0.91 ± 0.09	0.92 ± 0.05	0.96 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.02
SF3-gumbel	100	0.44 ± 0.05	0.90 ± 0.09	0.92 ± 0.03	0.97 ± 0.02	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF4-gumbel	10	0.71 ± 0.13	0.87 ± 0.16	0.77 ± 0.09	0.82 ± 0.07	0.91 ± 0.07	0.90 ± 0.07	0.89 ± 0.07	0.89 ± 0.08
SF4-gumbel	30	0.48 ± 0.08	0.88 ± 0.13	0.87 ± 0.07	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.96 ± 0.03	0.96 ± 0.03
SF4-gumbel	50	0.42 ± 0.07	0.89 ± 0.11	0.87 ± 0.05	0.95 ± 0.03	0.98 ± 0.02	0.98 ± 0.01	0.97 ± 0.02	0.97 ± 0.02
SF4-gumbel	100	0.39 ± 0.05	0.89 ± 0.10	0.88 ± 0.06	0.96 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.01	0.98 ± 0.01

1188
1189
11901191 Table 10: DAG learning performance (measured in true positive rate, the higher the better, best results
1192 in **bold**) of different algorithms on small scale (10-100 nodes) ER{2,3,4} graphs with different noise
1193 distributions. Our algorithm performs better than previous approaches.
1194

1195

Graphs	Nodes	MMPC	GES	NOTEARS	DAGMA	Order-1	Order-2	Order-3	Order-4
ER2-gauss	10	0.55 ± 0.05	0.46 ± 0.23	0.05 ± 0.07	0.02 ± 0.05	0.02 ± 0.04	0.02 ± 0.05	0.02 ± 0.06	0.03 ± 0.08
ER2-gauss	30	0.56 ± 0.03	0.46 ± 0.21	0.04 ± 0.06	0.02 ± 0.04	0.02 ± 0.04	0.02 ± 0.03	0.02 ± 0.03	0.02 ± 0.04
ER2-gauss	50	0.56 ± 0.03	0.47 ± 0.16	0.05 ± 0.06	0.01 ± 0.02	0.02 ± 0.02	0.01 ± 0.02	0.02 ± 0.02	0.02 ± 0.03
ER2-gauss	100	0.55 ± 0.02	0.44 ± 0.11	0.05 ± 0.05	0.01 ± 0.02	0.01 ± 0.02	0.01 ± 0.02	0.01 ± 0.02	0.02 ± 0.02
ER3-gauss	10	0.57 ± 0.05	0.60 ± 0.14	0.10 ± 0.08	0.06 ± 0.06	0.04 ± 0.05	0.04 ± 0.05	0.03 ± 0.05	0.04 ± 0.05
ER3-gauss	30	0.60 ± 0.04	0.70 ± 0.13	0.08 ± 0.08	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.04
ER3-gauss	50	0.59 ± 0.03	0.71 ± 0.10	0.10 ± 0.07	0.03 ± 0.03	0.03 ± 0.03	0.03 ± 0.03	0.03 ± 0.03	0.03 ± 0.03
ER3-gauss	100	0.58 ± 0.02	0.64 ± 0.12	0.11 ± 0.06	0.03 ± 0.03	0.03 ± 0.02	0.02 ± 0.02	0.02 ± 0.02	0.03 ± 0.03
ER4-gauss	10	0.53 ± 0.04	0.57 ± 0.13	0.09 ± 0.06	0.06 ± 0.05	0.05 ± 0.04	0.05 ± 0.04	0.04 ± 0.04	0.05 ± 0.04
ER4-gauss	30	0.61 ± 0.04	0.76 ± 0.07	0.14 ± 0.10	0.06 ± 0.05	0.05 ± 0.04	0.04 ± 0.04	0.05 ± 0.05	0.05 ± 0.06
ER4-gauss	50	0.62 ± 0.03	0.80 ± 0.05	0.15 ± 0.07	0.04 ± 0.04	0.04 ± 0.03	0.04 ± 0.03	0.04 ± 0.04	0.05 ± 0.04
ER4-gauss	100	0.61 ± 0.03	0.77 ± 0.06	0.15 ± 0.07	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04
ER2-exp	10	0.55 ± 0.05	0.48 ± 0.24	0.05 ± 0.07	0.03 ± 0.05	0.02 ± 0.04	0.02 ± 0.05	0.03 ± 0.06	0.04 ± 0.08
ER2-exp	30	0.56 ± 0.03	0.47 ± 0.21	0.05 ± 0.07	0.02 ± 0.04	0.02 ± 0.04	0.02 ± 0.04	0.01 ± 0.03	0.02 ± 0.04
ER2-exp	50	0.56 ± 0.03	0.49 ± 0.15	0.05 ± 0.06	0.01 ± 0.02	0.01 ± 0.02	0.01 ± 0.02	0.02 ± 0.02	0.02 ± 0.03
ER2-exp	100	0.55 ± 0.02	0.44 ± 0.12	0.05 ± 0.05	0.01 ± 0.02	0.01 ± 0.02	0.01 ± 0.02	0.02 ± 0.02	0.02 ± 0.02
ER3-exp	10	0.57 ± 0.06	0.60 ± 0.14	0.11 ± 0.08	0.06 ± 0.06	0.04 ± 0.05	0.04 ± 0.05	0.03 ± 0.04	0.04 ± 0.05
ER3-exp	30	0.60 ± 0.04	0.69 ± 0.13	0.08 ± 0.07	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.04
ER3-exp	50	0.59 ± 0.03	0.70 ± 0.12	0.10 ± 0.08	0.03 ± 0.04	0.03 ± 0.04	0.03 ± 0.03	0.03 ± 0.03	0.03 ± 0.03
ER3-exp	100	0.58 ± 0.02	0.64 ± 0.11	0.11 ± 0.06	0.03 ± 0.03	0.03 ± 0.02	0.03 ± 0.03	0.03 ± 0.02	0.03 ± 0.02
ER4-exp	10	0.53 ± 0.03	0.57 ± 0.12	0.09 ± 0.06	0.07 ± 0.05	0.05 ± 0.04	0.05 ± 0.04	0.04 ± 0.04	0.05 ± 0.04
ER4-exp	30	0.61 ± 0.04	0.76 ± 0.07	0.13 ± 0.09	0.06 ± 0.05	0.04 ± 0.04	0.05 ± 0.04	0.06 ± 0.05	0.05 ± 0.06
ER4-exp	50	0.62 ± 0.03	0.80 ± 0.05	0.15 ± 0.08	0.05 ± 0.04	0.04 ± 0.03	0.04 ± 0.03	0.04 ± 0.04	0.05 ± 0.04
ER4-exp	100	0.61 ± 0.03	0.77 ± 0.06	0.15 ± 0.07	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04
ER2-gumbel	10	0.55 ± 0.04	0.51 ± 0.23	0.03 ± 0.06	0.02 ± 0.05	0.02 ± 0.06	0.03 ± 0.06	0.03 ± 0.06	0.03 ± 0.08
ER2-gumbel	30	0.56 ± 0.03	0.49 ± 0.19	0.05 ± 0.07	0.02 ± 0.03	0.02 ± 0.04	0.02 ± 0.04	0.03 ± 0.04	0.03 ± 0.05
ER2-gumbel	50	0.56 ± 0.03	0.48 ± 0.16	0.06 ± 0.06	0.02 ± 0.03	0.02 ± 0.03	0.02 ± 0.03	0.03 ± 0.04	0.03 ± 0.04
ER2-gumbel	100	0.55 ± 0.02	0.44 ± 0.12	0.07 ± 0.06	0.01 ± 0.02	0.02 ± 0.02	0.02 ± 0.02	0.02 ± 0.02	0.03 ± 0.03
ER3-gumbel	10	0.57 ± 0.05	0.60 ± 0.15	0.07 ± 0.07	0.04 ± 0.05	0.03 ± 0.05	0.04 ± 0.06	0.03 ± 0.05	0.02 ± 0.04
ER3-gumbel	30	0.60 ± 0.04	0.70 ± 0.13	0.08 ± 0.08	0.02 ± 0.03	0.03 ± 0.04	0.03 ± 0.04	0.05 ± 0.05	0.04 ± 0.05
ER3-gumbel	50	0.59 ± 0.03	0.70 ± 0.12	0.10 ± 0.08	0.03 ± 0.03	0.04 ± 0.03	0.03 ± 0.03	0.04 ± 0.04	0.04 ± 0.04
ER3-gumbel	100	0.58 ± 0.02	0.64 ± 0.12	0.12 ± 0.06	0.03 ± 0.03	0.03 ± 0.03	0.03 ± 0.03	0.03 ± 0.03	0.04 ± 0.03
ER4-gumbel	10	0.53 ± 0.04	0.56 ± 0.13	0.08 ± 0.06	0.05 ± 0.04	0.03 ± 0.04	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04
ER4-gumbel	30	0.61 ± 0.04	0.75 ± 0.08	0.13 ± 0.09	0.04 ± 0.04	0.04 ± 0.04	0.04 ± 0.04	0.06 ± 0.06	0.06 ± 0.06
ER4-gumbel	50	0.61 ± 0.03	0.80 ± 0.05	0.15 ± 0.08	0.05 ± 0.05	0.04 ± 0.04	0.04 ± 0.03	0.05 ± 0.04	0.06 ± 0.04
ER4-gumbel	100	0.61 ± 0.03	0.77 ± 0.06	0.16 ± 0.07	0.04 ± 0.04	0.05 ± 0.05	0.04 ± 0.04	0.04 ± 0.04	0.05 ± 0.04
SF2-gauss	10	0.78 ± 0.11	0.91 ± 0.19	0.88 ± 0.09	0.91 ± 0.07	0.97 ± 0.04	0.97 ± 0.04	0.97 ± 0.04	0.96 ± 0.05
SF2-gauss	30	0.62 ± 0.09	0.95 ± 0.08	0.93 ± 0.05	0.95 ± 0.04	0.98 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.98 ± 0.02
SF2-gauss	50	0.57 ± 0.08	0.92 ± 0.11	0.93 ± 0.05	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF2-gauss	100	0.54 ± 0.06	0.92 ± 0.08	0.94 ± 0.04	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF3-gauss	10	0.73 ± 0.13	0.86 ± 0.19	0.75 ± 0.11	0.80 ± 0.10	0.91 ± 0.08	0.91 ± 0.08	0.91 ± 0.07	0.91 ± 0.07
SF3-gauss	30	0.53 ± 0.10	0.92 ± 0.09	0.88 ± 0.06	0.92 ± 0.04	0.97 ± 0.02	0.97 ± 0.03	0.97 ± 0.03	0.97 ± 0.03
SF3-gauss	50	0.47 ± 0.07	0.90 ± 0.10	0.89 ± 0.06	0.93 ± 0.03	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.02
SF3-gauss	100	0.43 ± 0.05	0.90 ± 0.09	0.88 ± 0.04	0.95 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF4-gauss	10	0.71 ± 0.13	0.89 ± 0.15	0.68 ± 0.09	0.73 ± 0.07	0.84 ± 0.07	0.85 ± 0.07	0.84 ± 0.07	0.84 ± 0.07
SF4-gauss	30	0.49 ± 0.08	0.87 ± 0.12	0.80 ± 0.08	0.87 ± 0.04	0.94 ± 0.03	0.95 ± 0.03	0.94 ± 0.03	0.94 ± 0.03
SF4-gauss	50	0.42 ± 0.07	0.88 ± 0.10	0.83 ± 0.05	0.91 ± 0.03	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02
SF4-gauss	100	0.39 ± 0.05	0.90 ± 0.09	0.85 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
SF2-exp	10	0.79 ± 0.11	0.91 ± 0.18	0.87 ± 0.09	0.91 ± 0.07	0.97 ± 0.04	0.97 ± 0.04	0.97 ± 0.04	0.96 ± 0.05
SF2-exp	30	0.62 ± 0.09	0.93 ± 0.11	0.93 ± 0.05	0.95 ± 0.04	0.98 ± 0.02	0.98 ± 0.02	0.99 ± 0.02	0.98 ± 0.02
SF2-exp	50	0.57 ± 0.08	0.92 ± 0.11	0.93 ± 0.05	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF2-exp	100	0.54 ± 0.06	0.91 ± 0.08	0.94 ± 0.04	0.97 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF3-exp	10	0.74 ± 0.13	0.90 ± 0.17	0.75 ± 0.11	0.81 ± 0.09	0.90 ± 0.07	0.91 ± 0.08	0.92 ± 0.07	0.90 ± 0.07
SF3-exp	30	0.53 ± 0.10	0.92 ± 0.10	0.88 ± 0.06	0.92 ± 0.04	0.97 ± 0.02	0.96 ± 0.03	0.97 ± 0.03	0.97 ± 0.03
SF3-exp	50	0.47 ± 0.07	0.91 ± 0.09	0.88 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.98 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
SF3-exp	100	0.43 ± 0.05	0.89 ± 0.10	0.89 ± 0.04	0.95 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01
SF4-exp	10	0.71 ± 0.13	0.88 ± 0.16	0.68 ± 0.09	0.72 ± 0.07	0.84 ± 0.08	0.84 ± 0.07	0.84 ± 0.08	0.84 ± 0.07
SF4-exp	30	0.49 ± 0.08	0.88 ± 0.12	0.81 ± 0.07	0.87 ± 0.05	0.94 ± 0.03	0.94 ± 0.03	0.94 ± 0.03	0.94 ± 0.04
SF4-exp	50	0.42 ± 0.07	0.89 ± 0.10	0.83 ± 0.05	0.90 ± 0.03	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02	0.96 ± 0.02
SF4-exp	100	0.39 ± 0.05	0.90 ± 0.09	0.85 ± 0.06	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.02	0.97 ± 0.01
SF2-gumbel	10	0.78 ± 0.11	0.91 ± 0.18	0.93 ± 0.08	0.96 ± 0.05	0.98 ± 0.04	0.98 ± 0.04	0.98 ± 0.04	0.98 ± 0.04
SF2-gumbel	30	0.62 ± 0.09	0.94 ± 0.08	0.95 ± 0.05	0.98 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.99 ± 0.02	0.99 ± 0.02
SF2-gumbel	50	0.57 ± 0.08	0.92 ± 0.10	0.95 ± 0.04	0.98 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF2-gumbel	100	0.53 ± 0.06	0.92 ± 0.08	0.95 ± 0.03	0.99 ± 0.02	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF3-gumbel	10	0.73 ± 0.13	0.89 ± 0.17	0.84 ± 0.10	0.89 ± 0.07	0.95 ± 0.06	0.94 ± 0.06	0.94 ± 0.06	0.95 ± 0.06
SF3-gumbel	30	0.54 ± 0.10	0.92 ± 0.10	0.92 ± 0.06	0.96 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.97 ± 0.02	0.97 ± 0.02
SF3-gumbel	50	0.48 ± 0.07	0.91 ± 0.09	0.92 ± 0.05	0.96 ± 0.02	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.01	0.98 ± 0.02
SF3-gumbel	100	0.44 ± 0.05	0.90 ± 0.09	0.92 ± 0.03	0.97 ± 0.02	0.98 ± 0.01	0.99 ± 0.01	0.99 ± 0.01	0.99 ± 0.01
SF4-gumbel	10	0.71 ± 0.13	0.87 ± 0.16	0.77 ± 0.09	0.82 ± 0.07	0.91 ± 0.07	0.90 ± 0.07	0.89 ± 0.07	0.89 ± 0.08
SF4-gumbel	30	0.48 ± 0.08	0.88 ± 0.13	0.87 ± 0.07	0.93 ± 0.03	0.97 ± 0.02	0.97 ± 0.02	0.96 ± 0.03	0.96 ± 0.03
SF4-gumbel	50	0.42 ± 0.07	0.89 ± 0.11	0.87 ± 0.05	0.95 ± 0.03	0.98 ± 0.02	0.98 ± 0.01	0.97 ± 0.02	0.97 ± 0.02
SF4-gumbel	100	0.39 ± 0.05	0.89 ± 0.10	0.88 ± 0.06	0.96 ± 0.02	0.98 ± 0.02	0.98 ± 0.02	0.98 ± 0.01	0.98 ± 0.01

1239
1240
1241

	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				
				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

```

1275
1276     def _h_grad(self, W, s, eps=1e-20):
1277         M_ = W * W / s
1278         Iw = self.Id - M_ # self.Id is identity matrix
1279         icnt = 1
1280         Inv = self.Id + M_
1281         while icnt < 2 * self.d:
1282             M_ = M_ @ M_
1283             Inv = Inv + Inv @ M_
1284             icnt *= 2
1285             if self.np.max(self.np.abs(M_)) < eps:
1286                 break
1287             if self.np.any(self.np.isnan(Inv)):
1288                 break
1289
1290             if self.np.any(self.np.isinf(Inv)):
1291                 return self.np.zeros_like(Inv)
1292
1293             if self.np.any(self.np.isnan(Inv)):
1294                 return self.np.zeros_like(Inv)
1295
1296         return Inv / s
1297
1298     def compute_h_grad(self, W, s):

```

```
1296
1297     M = self._h_grad(W, s)
1298     if self.np.any(self.np.isnan(M)) or self.np.linalg.norm(
1299         M @ (s * self.Id - W * W) - self.Id, ord='fro') >
1300         1e-6:
1301         if isinstance(W, cupy.ndarray):
1302             s, v = cupy.linalg.svd(W * W) # cupy does not
1303                 have a eig lib, thus use spectral norm as an
1304                 estimation
1305             cs = cupy.max(s) + 0.1 * self.h_order
1306         else:
1307             cs = np.max(np.abs(np.linalg.eigvals(
1308                 W * W))) + 0.1 * self.h_order
1309     else:
1310         cs = s
1311     return M, cs
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
```