

Efficient and Scalable Structure Learning for Bayesian Networks: Algorithms and Applications

Rong Zhu*, Andreas Pfadler*, Ziniu Wu*[†], Yuxing Han*, Xiaoke Yang*,
Feng Ye*, Zhenping Qian*, Jingren Zhou* and Bin Cui[‡]

*Alibaba Group, China, [†]Oxford University, UK, [‡]Peking University, China

*{red.zr, andreaswernerob, ziniu.wzn, yuxing.hyx, xiaoke.yxk, keven.yef,
zhengping.qzp, jingren.zhou}@alibaba-inc.com, [‡]bin.cui@pku.edu.cn

Abstract—Structure Learning for Bayesian network (BN) is an important problem with extensive research. It plays central roles in a wide variety of applications in Alibaba Group. However, existing structure learning algorithms suffer from considerable limitations in real-world applications due to their low efficiency and poor scalability. To resolve this, we propose a new structure learning algorithm LEAST, which comprehensively fulfills our business requirements as it attains high accuracy, efficiency and scalability at the same time. The core idea of LEAST is to formulate the structure learning into a continuous constrained optimization problem, with a novel differentiable constraint function measuring the acyclicity of the resulting graph. Unlike with existing work, our constraint function is built on the spectral radius of the graph and could be evaluated in near linear time w.r.t. the graph node size. Based on it, LEAST can be efficiently implemented with low storage overhead. According to our benchmark evaluation, LEAST runs 1–2 orders of magnitude faster than state-of-the-art method with comparable accuracy, and it is able to scale on BNs with up to hundreds of thousands of variables. In our production environment, LEAST is deployed and serves for more than 20 applications with thousands of executions per day. We describe a concrete scenario in a ticket booking service in Alibaba, where LEAST is applied to build a near real-time automatic anomaly detection and root error cause analysis system. We also show that LEAST unlocks the possibility of applying BN structure learning in new areas, such as large-scale gene expression data analysis and explainable recommendation system.

I. INTRODUCTION

Discovering, modeling and understanding *causal* mechanisms behind natural phenomena are fundamental tasks in numerous scientific disciplines, ranging from physics to economics, from biology to sociology. Bayesian network (BN) is prominent example of probabilistic graphical models and has been recognized as a powerful and versatile tool for modeling causality. Each BN takes the form of a directed acyclic graph (DAG), where each node corresponds to an observed or hidden variable and each edge defines the causal dependency between two variables. By further specifying the conditional probability distributions based on the causal structure, one eventually obtains a joint probability distribution for the model variables.

In real-world situations, running randomized controlled trials uncovering causal relationships for BN can be costly

and time consuming. Sometimes, it is even impossible due to ethical concerns. Thus, it has been of great interest to develop statistical methods to infer the structure of BN purely based on observed data. This problem, called *structure learning*, has become a research hot spot in machine learning (ML) community. Traditional applications include gene expression data analysis [18], [26], error identification [36], model factorization [9] and system optimization [34], to name but a few.

Application Scenarios. In Alibaba Group, a very large-scale company, structure learning for BN has been applied in a wide range of business scenarios, including but not limited to the areas of e-commerce, cloud computing and finance. In these areas, monitoring production systems is a representative application of BN structure learning. Traditional monitoring methods require a significant amount of human resources to analyze and understand log information compiled from a wide variety of internal and external interfaces. Thus, monitoring results heavily rely on personal experience, which may be inaccurate and incomplete. Meanwhile, due to complex data pipelines, business logic and technical infrastructure, staff feedback may not be obtained in time, greatly affecting the usability of some applications such as ticket booking and cloud security. For example, in an airline ticket booking application, operation staff often needs from several hours to days to find possible reasons causing booking errors.

To overcome these limitations, BN is a robust and interpretable tool for automatic anomaly detection and fast root error cause analysis in monitoring production systems. Due to the complexity and scale of our systems, BN needs to be generated in a data-driven fashion without relying too much on expert input. Therefore, it is very necessary to design a scalable and efficient structure learning method for BN to fulfill our business requirements.

Challenges. Learning the structure of a BN purely from data is a highly nontrivial endeavour, and it becomes much more difficult for our concrete business applications. Existing structure learning algorithms for BNs are based on discrete combinatorial optimization or numerical continuous optimization. Combinatorial optimization algorithms [1], [4], [5], [8], [10], [12], [13], [16], [20], [24], [28] can only produce accurate results for BNs with tens of nodes. Their accuracy significantly declines in the larger regimes. Recently, [17], [18], [37]–[39]

Rong Zhu and Andreas Pfadler contributed equally to this work, and Ziniu Wu is the corresponding author.

proposed a new paradigm by measuring the acyclicity of a graph in a numerical manner. In this paradigm, the structure learning problem is transformed into a constrained continuous optimization problem and can be solved using off-the-shelf gradient based methods. Therefore, this class of algorithms can be easily implemented in modern ML platforms and tend to produce accurate results. Yet, they suffer from efficiency and scalability issues. The time and space cost of computing the acyclicity metric are $O(d^3)$ and $O(d^2)$ for a d -node graph [38], [41], respectively, which only scales well to BNs with up to a thousand of nodes. As our application scenarios often contain tens of thousands to millions of variables, no existing method is applicable.

Our Contributions. In this paper, we tackle the above challenges by proposing a new structure learning algorithm called LEAST. It comprehensively fulfills our business requirements as it simultaneously attains high accuracy, high time efficiency and superior scalability. We leverage the advantages of continuous optimization algorithms, while designing a new way around their known drawbacks. Specifically, we find a novel and equivalent function to measure the acyclicity of the resulting graph based on its *spectral properties*. This new acyclicity constraint function is differentiable and much easier to compute. Based on it, LEAST can be efficiently implemented with low space overhead. Therefore, LEAST is able to learn structure of very large BN. We have deployed LEAST in the production environment of Alibaba. It now serves more than 20 business applications and is executed thousands of times per day.

Our main contributions can be summarized as follows:

- We propose LEAST, a scalable, efficient and accurate structure learning algorithm for BN. It is built upon a new acyclicity constraint, whose time and space computation cost is near linear w.r.t. the graph node size. (Sections 3 and 4)
- We conduct extensive experiments on benchmark datasets to verify the effectiveness, efficiency and scalability of LEAST. It attains a comparable result accuracy w.r.t. the state-of-the-art method while speeding up the computation by 1–2 orders of magnitude. Furthermore, it can scale to BNs with up to hundreds of thousands nodes. (Section 5)
- We demonstrate the usage of LEAST in different applications, including a production task in ticket booking system, a large-scale gene expression data analysis application and a case study in building explainable recommendation systems. We believe that this work represents a first step towards unlocking structure learning for a wide range of scenarios. (Section 6)

II. BACKGROUND

We first review some relevant concepts and background knowledge on BN and structure learning in this section.

Bayesian Networks and Structural Equation Models. We suppose that in our problem of interest we can observe a d -dimensional vector¹ $X \in \mathbb{R}^d$ of random variables. The random

vector X is assumed to be distributed according to some BN with DAG $G = (V, E)$, where V and E represent the set of nodes and edges, respectively. Thus, each node $i \in V$ exactly corresponds to a random variable X_i of X , and each edge $(i, j) \in E$ indicates some kind of causal relation from variable X_i to X_j . Let $X_{\text{pa}(i)}$ denote the set of parent random variables of X_i . Random variables in a BN model satisfy the *first order Markov* property. That is, each variable X_i is dependent on $X_{\text{pa}(i)}$, but independent from all other random variables conditioned on $X_{\text{pa}(i)}$. Therefore, one can decompose the high-dimensional joint probability density function (PDF) $p(X)$ of X into a product of compact local mass functions as $p(X) = \prod_{i=1}^d p(X_i | X_{\text{pa}(i)})$, where $p(X_i | X_{\text{pa}(i)}) = p(X_i)$ if X_i has no parents.

BNs admit an interesting interpretation: one may regard each variable X_i as a stochastic function of $X_{\text{pa}(i)}$, so that the causal relations encoded in a BN model naturally lead to a *structural equation model* (SEM) [14], [21]. For $i = 1, 2, \dots, d$, let n_i be a random noise variable and f_i be a measurable function. In an SEM, we set $X_i = f_i(X_{\text{pa}(i)}) + n_i$, to indicate the causal relation from parents $X_{\text{pa}(i)}$ to the random variable X_i . The functions $f_i(\cdot)$ are modeled as some linear or non-linear transformations of $X_{\text{pa}(i)}$ to which i.i.d. noise terms n_i are added. Once $f_i(\cdot)$ and n_i are suitably parameterized, an SEM may not only serve as a proxy to BN, but one may also *learn* a BN via estimation of SEM parameters.

In this paper, we consider the widely used linear SEM (LSEM) [30]. This means we set $X_i = w_i^T X$ where $w_i[j] = 0$ if X_j is not a parent of X_i , and noise terms n_i are not restricted to be Gaussian.

Structure Learning Problem. Given a sample matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d] \in \mathbb{R}^{n \times d}$ containing n i.i.d. observations of the random vector X assumed to follow a SEM model, the structure learning asks to output a BN G which encodes the causal relations between all X_i . For the LSEM case, let the matrix $\mathbf{W} = [w_1, w_2, \dots, w_d]$. Learning the structure of G is equivalent to finding the matrix \mathbf{W} where node i connects to node j in G iff $\mathbf{W}[i, j] \neq 0$. For clarity, we denote the graph induced by \mathbf{W} as $G(\mathbf{W})$. For a matrix \mathbf{W} , a decomposable loss function $L(\cdot)$, such as the least squares or the negative log-likelihood on all nodes, may be used to evaluate how well the model corresponding to \mathbf{W} fits the observed data. Hence, the structure learning problem is essentially solving the following program, which has been proven to be NP-Hard in [3]:

$$\begin{aligned} \arg \min_{\mathbf{W}} L(\mathbf{W}, \mathbf{X}) &= \frac{1}{n} \sum_{i=1}^d L(\mathbf{x}_i, w_i^T \mathbf{X}), \\ &\text{subject to } G(\mathbf{W}) \in \text{DAGs.} \end{aligned} \quad (1)$$

Existing Structure Learning Algorithms. From Eq. (1), the key challenge in BN structure learning is how to enforce the acyclicity constraint for $G(\mathbf{W})$. Historically this has been addressed using combinatorial optimization algorithms which directly explore the search space of all DAGs to optimize multiple kinds of loss scores [1], [5], [12], [16]. However, as

¹Unless otherwise statement, we assume all vectors to be column vectors.

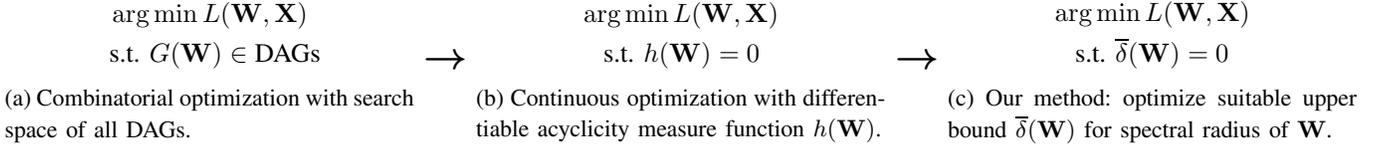


Fig. 1. Different optimization paradigms of structure learning algorithms for BNs.

the search space grows super-exponentially in the number of BN nodes, these methods can only scale to graphs with around ten nodes [23]. Later on, a large number of approximate combinatorial algorithms [4], [8], [10], [12], [13], [20], [24], [28] have been proposed to speed up by some heuristic pruning rules. They improve the efficiency while suffer a significant decline in accuracy [22], [25] in the large-scale regime. Besides, their scalability is still far from enough for many real-world applications [7], [18].

As a breakthrough, [39] proposed a novel method called NOTEARS to recast structure learning into a continuous optimization problem by measuring the acyclicity in a numerical manner. As shown in Fig. 1(b), $h(\mathbf{W})$ is a smooth non-negative function evaluating the closeness of graph $G(\mathbf{W})$ with DAG. $h(\mathbf{W}) = 0$ iff $G(\mathbf{W})$ is exactly a DAG. Based on this, the structure learning problem is transformed into a constrained continuous optimization problem and can be solved using off-the-shelf gradient based optimization methods. Algorithms in this class [17], [18], [37], [38] can generally handle a wide range of differentiable loss scores and tend to produce satisfying results. Yet, they suffer from efficiency and scalability issues since computing $h(\mathbf{W})$ requires $O(d^3)$ time and $O(d^2)$ space for a d -node graph [38], [41]. According to our tests, NOTEARS can process at most a thousand of nodes on a NVIDIA V100 GPU with 64GB RAM, requiring several hours until convergence.

Summary Existing structure learning algorithms for BNs suffer from considerable limitations w.r.t efficiency and scalability. None of them is applicable in our company’s business scenarios such as root error cause analysis and recommendation systems containing tens of thousands to millions of variables. To this end, we explore a new path to design efficient and scalable structure learning methods fulfilling our practical requirements. In the following Section 3, we propose a new acyclicity constraint, which establishes a new paradigm of structure learning for BNs. The detailed algorithm and its applications are then presented in Section 4 and 6, respectively.

III. NEW ACYCLICITY CONSTRAINT

In this section, we address the key challenge of structure learning for BNs by proposing a novel acyclicity constraint. We present our fundamental idea in Section III-A. Details on the formulation are given in Section III-B. Finally, in Section III-C we present the actual computation scheme for the constraint.

A. Preparations and Foundations

We first revisit the original constraint as presented in [38]. Based on this, we then explain the key points of our work.

Revisiting Existing Acyclicity Constraints. Let $\mathbf{S} = \mathbf{W} \circ \mathbf{W}$ where \circ is the Hadamard product. Then, G is a DAG iff

$$h(\mathbf{S}) = \text{Tr}(e^{\mathbf{S}}) - d = 0, \quad (2)$$

where $e^{\mathbf{S}}$ is the matrix exponential of \mathbf{S} .

We can regard \mathbf{S} as a non-negative adjacency matrix of the graph G , where each positive element $\mathbf{S}[i, j]$ indicates an edge (i, j) in G . For each $k \geq 1$ and every node i , $\mathbf{S}^k[i, i]$ is the sum of weights of all k -length cycles passing through node i in G . Therefore, there exist some cycles in G iff $\text{Tr}(\mathbf{S}^k) > 0$ for some $k \geq 1$. Since $e^{\mathbf{S}} = \sum_{i=0}^{\infty} \frac{\mathbf{S}^i}{i!}$ where $\mathbf{S}^0 = \mathbf{I}$, Eq. (2) certainly indicates that there exist no cycles in G . Later, [37] relaxes this constraint to

$$g(\mathbf{S}) = \text{Tr}((\mathbf{I} + \mathbf{S})^d) - d = \text{Tr}\left(\sum_{k=1}^d \binom{d}{k} \mathbf{S}^k\right) = 0. \quad (3)$$

For an acyclic graph G Eq. (3) holds since a simple cycle in G contains at most d nodes.

The acyclicity metrics $h(\mathbf{S})$ or $g(\mathbf{S})$ have two inherent drawbacks: 1) costly operations such as matrix exponential or matrix multiplication, whose complexity is $O(d^3)$; and 2) costly storage of the dense matrix $e^{\mathbf{S}}$ or $(\mathbf{I} + \mathbf{S})^d$ even though \mathbf{S} is often sparse. These two drawbacks fundamentally limit the efficiency and scalability of existing continuous optimization algorithms. To overcome the drawbacks, we need to design a new acyclicity metric for \mathbf{S} with *lightweight* computation and space overhead.

Fundamental Idea. We try to characterize the acyclicity of graph using the *spectral properties* of the matrix \mathbf{S} . Without loss of generality, let $\delta_1, \delta_2, \dots, \delta_d$ be the d eigenvalues of matrix \mathbf{S} . Since \mathbf{S} is non-negative, we have: 1) $\delta_i \geq 0$ for all $1 \leq i \leq d$; 2) $\text{Tr}(\mathbf{S}) = \sum_{i=1}^d \delta_i$; and 3) $\delta_1^k, \delta_2^k, \dots, \delta_d^k$ are the d eigenvalues of matrix \mathbf{S}^k for all $k \geq 2$. The absolute value of the largest eigenvalue δ is called the *spectral radius* of \mathbf{S} . By its definition, we obviously have $\sum_{k=1}^{\infty} \text{Tr}(\mathbf{S}^k) = 0$ iff $\delta = 0$. Therefore, the spectral radius can be used as a measure for acyclicity.

Prior work [18] has used δ for the acyclicity constraint. However, computing an exact or approximate spectral radius also requires $O(d^3)$ or $O(d^2)$ time, respectively. To this end, instead of using δ itself, we try to utilize a suitable proxy by deriving an upper bound $\bar{\delta}$ of δ , and then optimize Eq. (1) by asymptotically decreasing $\bar{\delta}$ to a very small value. When $\bar{\delta}$ is

small enough, δ will be close to 0. As shown in Fig. 1(c), this establishes a new paradigm of structure learning for BNs.

Requirements on $\bar{\delta}$. Obtaining an upper bound on the spectral radius is a longstanding mathematical problem and closely related to the well-known Perron-Frobenius theorem [2]. However, finding a suitable $\bar{\delta}$ is a non-trivial task, which should satisfy the following requirements:

- R1:** We can ensure acyclicity by using $\bar{\delta}$ as a proxy for δ during the optimization process, i.e., $\bar{\delta}$ is *consistent* to the exact δ .
- R2:** *Differentiability* w.r.t. \mathbf{S} so that $\bar{\delta}$ can be optimized with off-the-shelf gradient based methods.
- R3:** $\bar{\delta}$ and its gradient $\nabla_{\mathbf{S}}\bar{\delta}$ should be *time-efficient* to compute without involving costly matrix operations.
- R4:** $\bar{\delta}$ and its gradient $\nabla_{\mathbf{S}}\bar{\delta}$ should be *space-efficient* to compute without costly intermediate storage overhead.

B. Acyclicity Constraint Formulation

We formalize our proposed upper bound $\bar{\delta}$ in this subsection and show that it satisfies the above requirements. Given a matrix \mathbf{A} , let $r(\mathbf{A})$ and $c(\mathbf{A})$ be the vector of row sums and column sums of \mathbf{A} , respectively. Given a vector v and any real value α , let v^α be the resulting element-wise power vector. Our bound is derived in an iterative manner. Let $\mathbf{S}^{(0)} = \mathbf{S} = \mathbf{W} \circ \mathbf{W}$. For any $k \geq 0$ and $0 \leq \alpha \leq 1$, let $b^{(k)} = (r(\mathbf{S}^{(k)}))^\alpha \circ (c(\mathbf{S}^{(k)}))^{1-\alpha}$, $\mathbf{D}^{(k)} = \text{Diag}(b^{(k)})$ and

$$\mathbf{S}^{(k+1)} = (\mathbf{D}^{(k)})^{-1} \mathbf{S}^{(k)} \mathbf{D}^{(k)}, \quad (4)$$

where we set $(\mathbf{D}^{(k)})^{-1}[i, i] = 0$ if $\mathbf{D}^{(k)}[i, i] = 0$. We set the upper bound $\bar{\delta}^{(k)} = \sum_{i=1}^d b^{(k)}[i]$. The following lemma states the correctness of this upper bound, following [33]. Due to space limits, we omit all proofs in this version.

Lemma 1: For any non-negative matrix \mathbf{S} , $k \geq 0$ and $0 \leq \alpha \leq 1$, the spectral radius δ of matrix \mathbf{S} is no larger than $\bar{\delta}^{(k)}$.

In Eq. (4), we apply a diagonal matrix transformation on $\mathbf{S}^{(k)}$ and $\bar{\delta}^{(k)}$ gradually approaches the exact δ [33]. In our experiments, we find that setting k to a small number around 5 is enough to ensure the accuracy of the results. The factor α is a hyper-parameter balancing the effects of $r(\mathbf{S}^{(k)})$ and $c(\mathbf{S}^{(k)})$. We set it closer to 0 when values in $r(\mathbf{S}^{(k)})$ are much larger than those in $c(\mathbf{S}^{(k)})$ and vice versa so that the upper bound $\bar{\delta}^{(k)}$ will be smaller.

Consistency between $\bar{\delta}$ and $h(\mathbf{S})$ and $g(\mathbf{S})$. Next, we verify one by one that our upper bound satisfies all stated requirements. At first, the following lemma establishes the consistency between the upper bound $\bar{\delta}^{(k)}$ and the original acyclicity metrics $h(\mathbf{S})$ and $g(\mathbf{S})$.

Lemma 2: For any non-negative matrix \mathbf{S} , $k \geq 0$ and any value $\epsilon, \alpha \in (0, 1)$, if the upper bound $\bar{\delta}^{(k)} \leq \ln(\frac{\epsilon}{d} + 1)$, then $h(\mathbf{S}) \leq \epsilon$ holds; if the upper bound $\bar{\delta}^{(k)} \leq \frac{1}{\alpha} \log_d \frac{\epsilon}{d^2}$, then $g(\mathbf{S}) \leq \epsilon$ holds.

By Lemma 2, it follows that it is safe to use $\bar{\delta}^{(k)}$ as a proxy for $h(\mathbf{S})$ and $g(\mathbf{S})$ since they will also decrease to a very small value when we optimize $\bar{\delta}^{(k)}$. According to our

benchmark evaluation results in Section V, the correlation between $\bar{\delta}^{(k)}$ and $h(\mathbf{S})$ often exceeds 0.9. This indicates that R1 is satisfied. For the other requirements (R2–R4), we reserve the examination of them for the following subsection.

C. Efficient Computation of $\bar{\delta}$ and $\nabla_{\mathbf{S}}\bar{\delta}$

In this section, we introduce a time and space efficient way to compute the upper bound $\bar{\delta}^{(k)}$ and its gradient $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$.

Computing $\bar{\delta}^{(k)}$. Given a matrix \mathbf{A} and a column vector v , let $\mathbf{A} \circ v$ (or $v \circ \mathbf{A}$) and $\mathbf{A} \circ v^T$ (or $v^T \circ \mathbf{A}$) denote the resulting matrix by multiplying v onto each column of \mathbf{A} and v^T onto each row of \mathbf{A} , respectively. The diagonal transformation in Eq. (4) can then be equivalently written as

$$\mathbf{S}^{(k+1)} = (\mathbf{D}^{(k)})^{-1} \mathbf{S}^{(k)} \mathbf{D}^{(k)} = (b^{(k)})^{-1} \circ \mathbf{S}^{(k)} \circ (b^{(k)})^T. \quad (5)$$

Eq. (5) gives an explicit way to compute $\bar{\delta}^{(k)}$, which only requires scanning the non-zero elements in $\mathbf{S}^{(k)}$. Let s be the number of non-zero elements in matrix \mathbf{S} . Then, the time cost to obtain $\bar{\delta}^{(k)}$ is $O(ks)$. Since k is a small number and \mathbf{S} is often sparse for DAG, the time cost $O(ks)$ is close to $O(d)$. Meanwhile, the space cost to obtain $\bar{\delta}^{(k)}$ is at most $O(s)$. Therefore, it is time and space efficient to compute $\bar{\delta}$.

Computing $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$. We now turn our attention to computing $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$, which clearly exists according to Eq. (4) so R2 is satisfied. We will now manually apply backward differentiation to the iteration defining $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$ (sometimes referred to as the ‘‘adjoint’’ method outside of the ML literature). As a result, we obtain a recipe for iteratively computing $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$. This iterative method will allow us to implement the gradient such that $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$ remains *sparse* throughout the entire process.

To simplify notation, we denote the Hadamard product of two vectors $u^\alpha \circ v^{-\alpha}$ as $(\frac{u}{v})^\alpha$. By Eq. (5) and Lemma 1, following the chain rule, we obtain explicit formulae for $\nabla_{\mathbf{S}}\bar{\delta}^{(k)}$ as follows.

Lemma 3: For any non-negative matrix \mathbf{S} , $k \geq 0$ and $0 \leq \alpha \leq 1$, we always have

$$\nabla_{\mathbf{S}^{(k)}}\bar{\delta}^{(k)} = \nabla_{\mathbf{S}^{(k)}}b^{(k)} = x^{(k)} \circ \mathbf{J} + (y^{(k)})^T \circ \mathbf{J}, \quad (6)$$

where $x^{(k)} = \alpha \left(\frac{c(\mathbf{S}^{(k)})}{r(\mathbf{S}^{(k)})} \right)^{1-\alpha}$, $y^{(k)} = (1 - \alpha) \left(\frac{r(\mathbf{S}^{(k)})}{c(\mathbf{S}^{(k)})} \right)^\alpha$, and $\mathbf{J} \in \mathbb{R}^{d \times d}$ is a matrix with all entries equal to 1.

Lemma 4: For any non-negative matrix \mathbf{S} , $k \geq 1$ and $0 \leq \alpha \leq 1$, given any $1 \leq j \leq k$, suppose that we already have $\nabla_{\mathbf{S}^{(j)}}\bar{\delta}^{(k)}$, let

$$z^{(j-1)} = - \frac{r(\nabla_{\mathbf{S}^{(j)}}\bar{\delta}^{(k)} \circ \mathbf{S}^{(j-1)} \circ (b^{(j-1)})^T)}{(b^{(j-1)})^2} + c((b^{(j-1)})^{-1} \circ \nabla_{\mathbf{S}^{(j)}}\bar{\delta}^{(k)} \circ \mathbf{S}^{(j-1)}). \quad (7)$$

Then, we have

$$\nabla_{\mathbf{S}^{(j-1)}}\bar{\delta}^{(k)} = (b^{(j-1)})^{-1} \circ \nabla_{\mathbf{S}^{(j)}}\bar{\delta}^{(k)} \circ (b^{(j-1)})^T + x^{(j-1)} \circ z^{(j-1)} \circ \mathbf{J} + (y^{(j-1)})^T \circ (z^{(j-1)})^T \circ \mathbf{J}, \quad (8)$$

Procedure FORWARD(\mathbf{W}, k, α)

```

1:  $\mathbf{S}^0 = \mathbf{W}^2$ 
2: for  $j = 0$  to  $k$  do
3:    $b^{(j)} \leftarrow (r(\mathbf{S}^{(j)}))^\alpha \circ (c(\mathbf{S}^{(j)}))^{1-\alpha}$ 
4:   if  $j \leq k-1$  then
5:     compute  $\mathbf{S}^{(j+1)}$  by Eq. (5)
6:   end if
7: end for
8: return  $\bar{\delta}^{(k)} \leftarrow \sum_{i=1}^d b^{(k)}[i]$ 

```

Procedure BACKWARD(\mathbf{W}, k, α)

```

1:  $\mathbf{M} \leftarrow \mathbf{W}^0$ 
2: compute  $x^{(k)}$  and  $y^{(k)}$  by Lemma 3
3: compute  $\nabla'_{\mathbf{S}^{(k)}} \bar{\delta}^{(k)}$  by Lemma 5
4: for  $j = k$  to  $1$  do
5:   compute  $x^{(j-1)}$  and  $y^{(j-1)}$  by Lemma 3
6:   compute  $z^{(j-1)}$  by Eq. (7)
7:   compute  $\nabla'_{\mathbf{S}^{(j-1)}} \bar{\delta}^{(k)}$  by Eq. (9)
8: end for
9: return  $\nabla_{\mathbf{W}} \bar{\delta}^{(k)} \leftarrow 2\nabla'_{\mathbf{S}^{(0)}} \bar{\delta}^{(k)} \circ \mathbf{W}$ 

```

Fig. 2. Procedures for computing $\bar{\delta}^{(k)}$ and $\nabla_{\mathbf{W}} \bar{\delta}^{(k)}$.

where $x^{(j-1)}$ and $y^{(j-1)}$ have the same meaning as Lemma 3.

Now, if we directly compute $\nabla_{\mathbf{S}} \bar{\delta}^{(k)}$ following the above lemmas, the resulting algorithm would not be space efficient since $\nabla_{\mathbf{S}^{(k)}} \bar{\delta}^{(k)}$ would be a dense matrix for all $0 \leq j \leq k$. Nevertheless, since the final objective is to compute $\nabla_{\mathbf{W}} \bar{\delta}^{(k)} = 2\nabla_{\mathbf{S}} \bar{\delta}^{(k)} \circ \mathbf{W}$, it is only necessary to compute the gradient of non-zero elements in \mathbf{W} and \mathbf{S} . Interestingly, we find that it is also safe to do this “masking” in advance. The correctness is guaranteed by the following lemma.

Lemma 5: For any matrix \mathbf{W} , $k \geq 1$ and $0 \leq \alpha \leq 1$, let $\mathbf{M} \in \mathbb{R}^{d \times d}$ be such that $\mathbf{M}[i, j] = 1$ when $\mathbf{W}[i, j] \neq 0$ and $\mathbf{M}[i, j] = 0$ otherwise. Let $\nabla'_{\mathbf{S}^{(k)}} \bar{\delta}^{(k)} = x^{(k)} \circ \mathbf{M} + (y^{(k)})^T \circ \mathbf{M}$. For all $1 \leq j \leq k-1$, let

$$\begin{aligned} \nabla'_{\mathbf{S}^{(j-1)}} \bar{\delta}^{(k)} &= (b^{(j-1)})^{-1} \circ \nabla'_{\mathbf{S}^{(j)}} \bar{\delta}^{(k)} \circ (b^{(j-1)})^T \\ &+ x^{(j-1)} \circ z^{(j-1)} \circ \mathbf{M} + (y^{(j-1)})^T \circ (z^{(j-1)})^T \circ \mathbf{M}, \end{aligned} \quad (9)$$

where $x^{(j-1)}$ and $y^{(j-1)}$ are defined in Lemma 3, and $z^{(j-1)}$ is defined in Eq. (7). Then, we always have

$$\nabla_{\mathbf{W}} \bar{\delta}^{(k)} = 2\nabla'_{\mathbf{S}} \bar{\delta}^{(k)} \circ \mathbf{W}. \quad (10)$$

According to Lemma 5, we can obtain $\nabla_{\mathbf{W}} \bar{\delta}^{(k)}$ in a space efficient way by properly making use of the sparsity structure of \mathbf{W} and \mathbf{S} . For each $1 \leq j \leq k-1$, the time cost to compute $x^{(j)}$, $y^{(j)}$ and $h^{(j)}$ are all $O(s)$. Therefore, the time cost to compute $\nabla_{\mathbf{W}} \bar{\delta}^{(k)}$ is also $O(ks)$, which is close to $O(d)$ for sparse \mathbf{W} .

Procedures Description. To summarize, we describe the whole process to compute $\bar{\delta}^{(k)}$ and $\nabla_{\mathbf{W}} \bar{\delta}^{(k)}$ in Fig. 2. Note that computing both $\bar{\delta}^{(k)}$ and $\nabla_{\mathbf{W}} \bar{\delta}^{(k)}$ will cost at most $O(ks)$ time and $O(s)$ space, near linear $O(d)$ cost for DAGs. Therefore, the requirements R3 and R4 are satisfied.

Algorithm LEAST($X, \zeta, \lambda, \epsilon, k, \alpha, B, \theta, T_o, T_i$)

```

1:  $\rho \leftarrow 1, \eta \leftarrow 1$ 
2: repeat
3:    $\mathbf{W}^*, \bar{\delta}(\mathbf{W}^*) \leftarrow \text{INNER}(X, \zeta, \lambda, \rho, \eta, k, \alpha, B, \theta, T_i)$ 
4:    $\eta \leftarrow \eta + \rho \bar{\delta}(\mathbf{W}^*)$ 
5:   enlarge  $p$  by a small factor
6: until  $\bar{\delta}(\mathbf{W}^*) \leq \epsilon$  or running  $T_o$  iterations
7: return  $\mathbf{W}^*$ 

```

Procedure INNER($X, \zeta, \lambda, \rho, \eta, k, \alpha, \theta, T$)

```

1: randomly initialize  $\mathbf{W}$  as a sparse matrix with density  $\zeta$  using Glorot uniform initialization
2: repeat
3:    $\bar{\delta}(\mathbf{W}) \leftarrow \text{FORWARD}(\mathbf{W}, k, \alpha)$ 
4:    $\nabla_{\mathbf{W}} \bar{\delta}(\mathbf{W}) \leftarrow \text{BACKWARD}(\mathbf{W}, k, \alpha)$ 
5:   randomly fetch a sample  $\mathbf{X}_B$  from  $X$  with  $B$  samples
6:    $\ell(\mathbf{W}) \leftarrow L(\mathbf{W}, \mathbf{X}_B) + \frac{\rho}{2} \bar{\delta}(\mathbf{W})^2 + \eta \bar{\delta}(\mathbf{W})$ 
7:    $\nabla_{\mathbf{W}} \ell(\mathbf{W}) \leftarrow \nabla_{\mathbf{W}} L(\mathbf{W}, \mathbf{X}_B) + (\rho + \bar{\delta}(\mathbf{W})) \nabla_{\mathbf{W}} \bar{\delta}(\mathbf{W})$ 
8:   update  $\mathbf{W}$  with some optimizer using  $\nabla_{\mathbf{W}} \ell(\mathbf{W})$  % e.g. Adam

9:   filter all elements in  $\mathbf{W}$  whose absolute value is less than  $\theta$ 
10: until  $\ell(\mathbf{W})$  converges or running  $T$  iterations
11: return  $\mathbf{W}$  and  $\bar{\delta}(\mathbf{W})$ 

```

Fig. 3. Algorithm for large-scale BN structure learning.

IV. STRUCTURE LEARNING ALGORITHM

In this section, we propose the structure learning algorithm built on top of the new acyclicity constraint introduced in the previous section. Our algorithm is called LEAST, which represents a Large-scale, Efficient and Accurate Structure learning method for BNs. First, we present the details of the algorithm in our exemplary LSEM case. Then, we shortly discuss its implementation details.

LEAST Algorithm. Given a matrix \mathbf{A} , let $\|\mathbf{A}\|_p$ denote the p -norm of \mathbf{A} for all $p \geq 1$. Following [38], we set the loss function in Eq. (1) as $L(\mathbf{W}, \mathbf{X}) = \frac{1}{n} \|\mathbf{X} - \mathbf{X}\mathbf{W}\|_2^2 + \lambda \|\mathbf{W}\|_1$, that is the least squares with L_1 -regularization. We denote $\bar{\delta}(\mathbf{W})$ as the spectral radius upper bound of the matrix \mathbf{W} . To solve this optimization problem, we use the augmented Lagrangian [19] method with some adjustments. Concretely, we solve the following program augmented from Eq. (1) with a quadratic penalty:

$$\begin{aligned} \arg \min_{\mathbf{W}} \quad & \frac{1}{n} \|\mathbf{X} - \mathbf{X}\mathbf{W}\|_2^2 + \lambda \|\mathbf{W}\|_1 + \frac{\rho}{2} \bar{\delta}(\mathbf{W})^2, \\ \text{subject to} \quad & \bar{\delta}(\mathbf{W}) = 0, \end{aligned} \quad (11)$$

where $\rho > 0$ is a penalty parameter. Eq. (11) is then transformed into an unconstrained program. Specifically, let η denote the Lagrangian multiplier, and $\ell(\mathbf{W}) = L(\mathbf{W}, \mathbf{X}) + \frac{\rho}{2} \bar{\delta}(\mathbf{W})^2 + \eta \bar{\delta}(\mathbf{W})$ represent the unconstrained objective function. Solving Eq. (11) then consists of two iterative steps:

1) find the optimal value \mathbf{W}^* minimizing the function $\ell(\mathbf{W})$;

2) increase the multiplier η to $\eta + \rho \bar{\delta}(\mathbf{W})$ and optionally increase the penalty value ρ to improve convergence.

With the growth of ρ and η , since $\bar{\delta}(\mathbf{W}) \geq 0$, minimizing $\ell(\mathbf{W})$ forces $\bar{\delta}(\mathbf{W})$ to decrease to near 0. LEAST follows

this generic process. In each iteration, we use the Inner procedure to optimize $\ell(\mathbf{W})$ using a first gradient based method (line 3), such as Adam [15], and then enlarge ρ and η (lines 4–5). Finally, we return \mathbf{W} when $\bar{\delta}(\mathbf{W})$ is below a small tolerance value ϵ . In each iteration, we initialize \mathbf{W} as a random sparse matrix with density ζ using Glorot uniform initialization (line 1). To improve efficiency and scalability, we also incorporate the two techniques into LEAST:

1) **Batching:** in line 5 of INNER, we randomly fetch a batch \mathbf{X}_B from \mathbf{X} to optimize $\ell(\mathbf{W})$ instead of using \mathbf{X} .

2) **Thresholding:** in line 9 of the Inner procedure, we filter elements in \mathbf{W} with absolute value below a small threshold θ . A small numerical value in \mathbf{W} indicates a weak correlation between the two nodes, so filtering out them may help ruling out false cycle-inducing edges in advance. It has been shown that this helps to decrease the number of false positive edges [35], [38], [40]. Moreover, removing these elements makes \mathbf{W} remain sparse throughout the optimization process, thus ensuring overall computational efficiency.

Complexity Analysis. We now analyze the computational complexity of each round in LEAST. Computing $L(\mathbf{W}, X_B)$ and its gradient $\nabla_{\mathbf{W}}L(\mathbf{W}, X_B)$ costs $O(Bsd)$ time and $O(s + Bd)$ space, where s is the number of non-zero elements in \mathbf{W} . As a result, the time cost for computing the acyclicity constraint $O(s) \ll O(Bsd)$. Thus, overall training time will be dominated by the first term in Eq. (11). Meanwhile, since \mathbf{W} remains sparse and \mathbf{X}_B is a small matrix, they can easily fit into memory.

Implementation Details. We implement two versions of our LEAST algorithm in Python:

1) LEAST-TF is a Tensorflow-based implementation using dense tensors only. We rely on Tensorflow’s automatic differentiation capabilities to compute gradients.

2) LEAST-SP based on the SciPy library, where we have implemented LEAST using sparse data structures. We use the CSR format to represent the sparse matrices. Gradients for the constraint function $\bar{\delta}$ are implemented according to the description in Section III. We use the Adam [15] optimizer, since it exhibits fast convergence and does not generate dense matrices during the computation process.

The two implementations are used in different scenarios. We generally use LEAST-TF for those cases where a dense \mathbf{W} may fit entirely into the (GPU) memory. This is particularly relevant for those applications which have strict requirements on training speed, such as real-time monitoring systems presented in Section VI-A. We also use LEAST-TF to validate our method with respect to our accuracy and efficiency claims in the following Section V.

For use cases where data contains hundreds of thousands variables and a high degree of sparsity, we use LEAST-SP. By leveraging sparse data structures, LEAST-SP is able to deal with cases where the dense representation \mathbf{W} does not fit into GPU or main memory anymore. It is thus also used for our scalability testing experiments in the next section.

We have deployed both LEAST-TF and LEAST-SP in the production environment of our company. They are part

of an internal library supporting different applications in Alibaba. Until now, they have been used in more than 20 business scenarios including production system monitoring, recommendation systems, health care data analysis, money laundering detection, could security and etc. It is executed thousands of times per day. We also plan to release an open-source implementation of LEAST to the research community.

V. EVALUATION RESULTS

We conduct extensive experiments to evaluate the accuracy, efficiency and scalability of our proposed method. The results are reported in this section.

Algorithms. We compare our LEAST algorithm with NOTEARS, the state-of-the-art structure learning algorithm for BN proposed in [38]. It provides the most straightforward way to establish the correctness of our proposed method, while allowing for a fair comparison w.r.t. computational efficiency. We refrain from comparing with traditional combinatorial optimization algorithms such as [22], [31], [32]. The reason is that we obtain results consistent with [38], which already shows that NOTEARS generally outperforms them.

Parameter Settings. In our LEAST method, for the upper bound $\bar{\delta}(\mathbf{W})$, we set $k = 5$ and the balancing factor $\alpha = 0.9$. For the optimizer in the INNER procedure, we use Adam [15] and set its learning rate to 0.01. We set the initialization density $\zeta = 10^{-4}$. We furthermore ensure that in LEAST-SP Adam is operating on sparse matrices only. The maximum outer and inner iteration numbers are set to 1,000 and 200 for all algorithms, respectively. Remaining parameters are tuned individually for each experiment.

For NOTEARS we use the Tensorflow implementation provided in [18]. Note that, it seems hardly possible to implement NOTEARS purely using sparse matrices, as some steps would always involve completely dense matrices.

Environment. All of our experiments are run on a Linux server with a Intel Xeon Platinum 2.5 GHz CPU with 96 cores and 512GB DDR4 main memory. For experiments on GPUs we use a NVIDIA Tesla V100 SXM2 GPU with 64GB GPU memory. We run each experiment multiple times and report the mean value of each evaluation metric.

In the following, Section V-A reports the evaluation results on artificial benchmark datasets in terms of accuracy and efficiency. Section V-B examines the scalability of our method on large-scale real-world datasets.

A. Artificial Datasets

We use the graph generation code from [38] to produce benchmark artificial data. It generates a random graph topology of G following two models, Erdős-Rényi (ER) or scale-free (SF), and then assigns each edge a uniformly random weight to obtain the adjacency matrix \mathbf{W} . The sample matrix \mathbf{X} is then generated according to LSEM with three kinds of additive noise: Gaussian (GS), Exponential (EX), and Gumbel (GB). We vary the node size d , the sample size n and the average degree of nodes. Following [38], we set the average

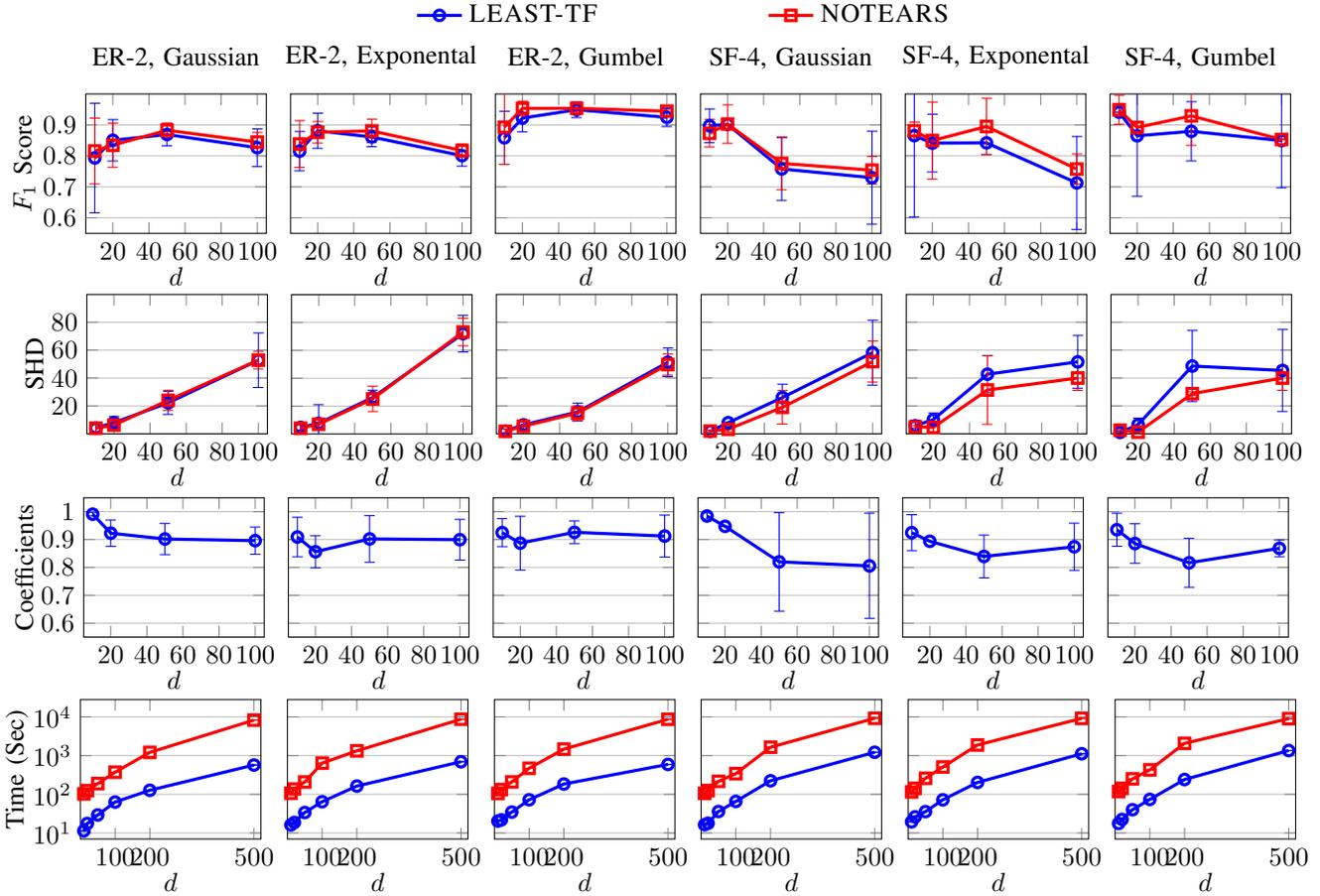


Fig. 4. Evaluation results on artificial benchmark data.

node degree to 2 for ER and 4 for SF graphs. For fairness, we compare LEAST-TF, the Tensorflow implementation of LEAST with NOTEARS. We also slightly modify the termination condition of LEAST. At the end of each outer loop, we also compute the value of $h(\mathbf{W})$ and terminate when $h(\mathbf{W})$ is smaller than the tolerance value ϵ . In this way, we ensure that their convergence is tested using the same termination condition. For remaining parameters, we set the batch size B equal to the sample size n , the filtering threshold $\theta = 0$ and the regularization penalty factor $\lambda = 0.5$. All experiments in this subsection are run on the CPU.

Result Accuracy. We evaluate by comparing results with the original ground truth graph G used to generate random samples. Following [38], after optimizing the result matrix \mathbf{W} to a small tolerance value ϵ , we filter it using a small threshold τ to obtain \mathbf{W}' , and then compare $G(\mathbf{W}')$ with G . We apply a grid search for the two hyper-parameters $\epsilon \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}\}$ and $\tau \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$, and report the result of the best case. We vary $d = 10, 20, 50, 100$ and set $n = 10d$. The results for the two graph models with three types of noise are shown in Fig. 4.

The first two lines of Fig. 4 report the results in terms of F_1 -score and the Structural Hamming Distance (SHD), respectively. The third line reports the detailed Pearson correlation coefficients of $\bar{\delta}(\mathbf{W})$ and $h(\mathbf{W})$ recorded during

the computation process of LEAST method. The error bar indicates the standard deviation of each value. We have the following observations:

1) For LEAST we obtain F_1 -scores which are larger than 0.8 in almost all cases. Meanwhile, the results for LEAST are very close to those of NOTEARS in terms of the F_1 -score and SHD. The difference appears to be negligible in most cases.

2) Our upper bound based acyclicity measure $\bar{\delta}(\mathbf{W})$ is highly correlated with the original metric $h(\mathbf{W})$, with correlation coefficients larger than 0.8 in all cases and larger than 0.9 in most cases. This verification indicates that $\bar{\delta}(\mathbf{W})$ is consistent with $h(\mathbf{W})$ and a valid acyclicity measure.

3) Results of LEAST appear to have higher variance than those of NOTEARS. This is likely due to small numerical instabilities in the iterative computation process for $\bar{\delta}(\mathbf{W})$.

4) the observed difference in variance between LEAST and NOTEARS is larger on dense SF-4 graphs than sparse ER-2 graphs, also more noticeable for large $d \geq 50$ than small d . Interestingly, as soon as the variance increases, correlation coefficients decrease at the same time. This could be attributed to the number of non-zero elements in \mathbf{W} , which increases w.r.t. d and the graph density.

Overall, this set of experiments shows that, for the above graphs, our proposed acyclicity metric $\bar{\delta}(\mathbf{W})$ produces results consistent with the original metric $h(\mathbf{W})$. Yet, we note that the

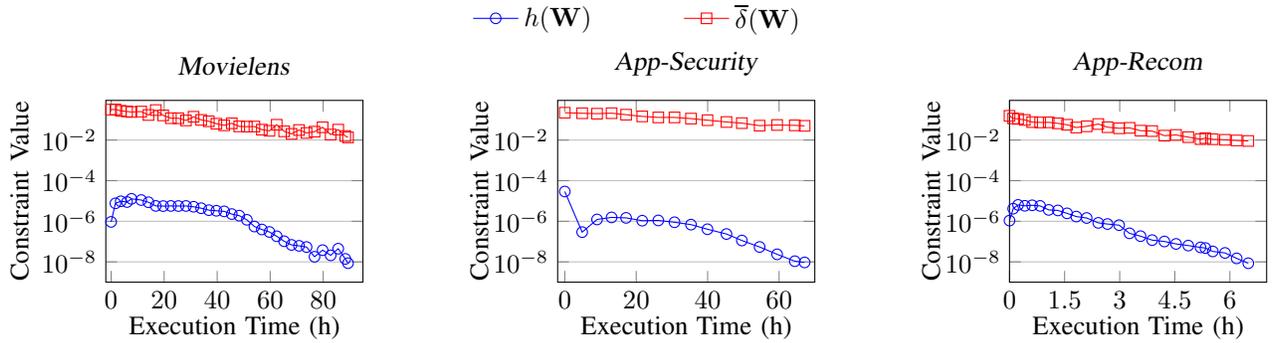


Fig. 5. Scalability test of LEAST algorithm.

increased computational efficiency (see the following results) comes at the price of an increased variance.

Time Efficiency. The evaluation results are reported in the fourth row of Fig. 4, which represents the execution time by fixing $\epsilon = 10^{-4}$ and sampling size $n = 10d$. We find that:

1) LEAST runs faster than NOTEARS in all tested cases. The speed up ratio ranges from 5 to 15. This is mainly due to the fact that the time complexity of computing the acyclicity constraint in the case of LEAST is much lower than that in NOTEARS (near $O(d)$ vs. $O(d^3)$).

2) the speed up effects of LEAST are similar for different graph models and noise types. This is because the time complexity of LEAST only depends on d , but is independent of all other factors.

3) the speed up effect is more obvious when d gets larger. The speed up ratio is up to 10, 9.5 and 14.7 when $d = 100$, 200 and 500, respectively. This is because the time cost for evaluating the acyclicity constraint grows almost linearly with d in LEAST while cubic with d in NOTEARS.

Overall, this set of experiments indicates that LEAST generally outperforms NOTEARS in terms of computational efficiency, especially for larger graphs.

Summary. On artificial benchmark data, our LEAST algorithm attains comparable accuracy w.r.t. the state-of-the-art NOTEARS method while improves the time efficiency by 1–2 orders of magnitude.

B. Real-World Datasets

We examine the scalability of our method on large-scale real-world datasets. Their properties of are summarized in Table III. It should be noted that it is difficult to find suitable, publicly available benchmark datasets which involve more than thousands of variables. This is why we also present some results based on non-public, internal datasets.

- *MovieLens* is based on the well-known MovieLens20M [11] recommendation dataset. We regard each movie as a node and each user’s entire rating history as one sample. The rating values range from 0 to 5 with an average rating of about 3.5. We construct the data matrix \mathbf{X} in a standard way as follows: Let r_{ij} be the rating of user i for movie j . Then, we set $\mathbf{X}[i, j] = r_{ij} - \sum_{j=1}^{n_i} r_{ij}/n_i$, where n_i is the number of available ratings for user i , so that we subtract

TABLE I
PROPERTIES OF REAL-WORLD LARGE-SCALE DATASETS.

Dataset Name	# of Nodes	# of Samples
<i>MovieLens</i>	27, 278	138, 493
<i>App-Security</i>	91, 850	1, 000, 000
<i>App-Recom</i>	159, 008	584, 871

each user’s mean rating from their ratings. We present a case study using this dataset later in Section VI-C.

- *App-Security* and *App-Recom* are two datasets extracted from real application scenarios at Alibaba. They are used in the context of cloud security and recommendation systems, respectively. Due to the sensitive nature of these businesses, we cannot provide more details beyond their size.

We set the batch size $B = 1,000$, the filtering threshold $\theta = 10^{-3}$ and the tolerance value $\epsilon = 10^{-8}$ and run LEAST-SP, the sparse matrices based implementation of LEAST on them. Note that, NOTEARS is unable to scale to these datasets. Fig. 5 illustrates how the constraint value $\bar{\delta}(\mathbf{W})$ and $h(\mathbf{W})$ varies with regard to the execution time. We find that:

1) In all cases, when optimizing the constraint $\bar{\delta}(\mathbf{W})$, the value of $h(\mathbf{W})$ also decreases accordingly and converges to a very small level. This verifies the effectiveness of our proposed acyclicity constraint on large graphs.

2) Our proposed algorithm can scale to very large graphs while taking a reasonable time. It takes 89.4, 67.2 and 6.5 hours on the datasets *MovieLens*, *App-Security* and *App-Recom*, respectively. So far, to the best of our knowledge, no existing continuous optimization structure learning algorithm can process SLP with more than 10^4 nodes, whereas our proposed method successfully scales to at least 10^5 nodes.

VI. APPLICATIONS

In this section, we demonstrate several different types of applications to show the effectiveness of our LEAST method. Section VI-A describes how LEAST is used in the production system of ticket booking service. Section VI-B describes the usage of LEAST for large-scale gene expression data analysis. Section VI-C presents a case study for applying LEAST in the context explainable recommendation system.

A. Monitoring the Flight Ticket Booking Business

Scenario Description. Alibaba owns Fliggy (Fei Zhu), a flight ticket booking service which serves several millions of

customers every day. It runs in a distributed system of considerable size and complexity, both growing with each passing day. The system is connected to wide range of interfaces to 1) a variety of airlines; 2) large intermediary booking systems such as Amadeus or Travelsky; 3) different smaller travel agents; 4) other internal systems in Alibaba and etc. We use the term fare source to reflect through which channel a ticket is booked. The booking process of each flight ticket consists of four essential steps: 1) query and confirm seat availability; 2) query and confirm price; 3) reserve ticket; and 4) payment and final confirmation. Each step involved one or several API requests of some interfaces stated above. All this information is logged for monitoring.

Due to many unpredictable reasons, e.g., an update or maintenance of an airline’s booking system, the booking process may fail at any step. To maintain a low rate of booking errors and ensure customers’ experience, an automatic monitoring system has been implemented. Whenever there is a sudden increase in booking errors of each step, log data is inspected to identify possible root causes. Once found, operation staff can then manually intervene.

In this situation, the monitoring system must fulfill two goals: 1) It should detect issues as *fast* as possible in order to reduce the possibility of money loss due to unsuccessful booking. 2) It should clearly find the actual *root* cause of a problem, since a failure originating at a single point may easily influence many other factors in the booking system.

Technical Details. Our LEAST algorithm plays a key role in this monitoring system. It works in the following manner:

1) For every half hour interval, we collect log data T from a window of the latest 24 hours and learn a BN network G using LEAST. The nodes in G include the four “error-type nodes” corresponding to the above four booking steps and information on airlines, fare sources and airports in the ticket booking system. Fig. 6 illustrates an example of the graph learned by the ticket booking data with partial nodes.

2) For each node X of the four error types, we inspect all paths P whose destination is X . That is, we follow the incoming links of X until we reach a node with no parents. Each P stands for a possible reason causing X . To detect whether P is a random coincidence or not, we count the number of occurrences of P in the log data T and T' , i.e., the previous window of log data, and perform a statistical test to derive a p-value. Depending on a threshold for the p-value, the entire path P may be reported as an anomaly, with the tail of P likely pinpointing the root cause for the problem.

In general, we have observed that learned BN structure (and weights on edges) remain unchanged across periods, as long as there are no observed anomalies leading to error rate increases. However, in other cases, we have found that relatively large increases in booking error rates are often accompanied by the appearance of new links pointing to the error nodes.

Effects and Results. Owing to the high efficiency of our LEAST method, the analysis task could be finished in 2–3 minutes of each run. Therefore, the operation staff can

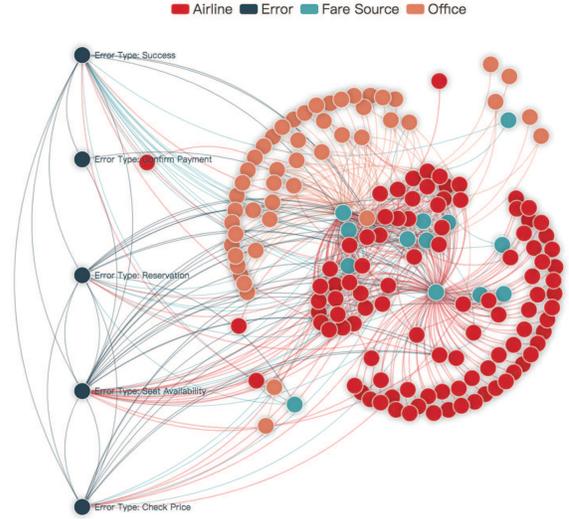


Fig. 6. Example of a graph learned from Fliggy booking data with partial nodes.

take immediate actions to resolve the problems, rather than waiting for several hours or days for domain expert analysis and feedback as before. After deploying our LEAST-based monitoring system, the *first-attempt success booking rate*, a metric evaluating the percentage of all tickets having been successfully booked upon the first try, has been improved.

To evaluate our LEAST-based monitoring system, we have compared the output of LEAST for a period of several weeks with results obtained manually from domain experts. Fig. 7 reports the results of this evaluation. Clearly, only 3% of the reported cases are found to be false positives but 97% of them are true positives. Among them, 42% of the reported problems occur relatively frequently due to problems with some external systems. 3%, 3% and 10% of the problems are correctly associated to issues caused by airlines, intermediary interfaces and travel agents, respectively, which has been verified manually. The remaining 39% were related to actual problems, the root causes of which remained unknown even after manual inspection and analysis, which can be due to a variety of highly unpredictable events (such as cabin or route adjustments from airlines and sudden weather changes).

Table II reports some examples of observed anomalies with explainable events. Some of these problems were ultimately caused by technical issues with an airline or intermediary interface. Other issues listed in Table II were caused by effects of the COVID-19 pandemic, such as sudden travel restrictions. While these events are not directly observable from the log data, our method was able to successfully detect their effects in our business.

Summary. Overall, our LEAST-based anomaly detection and root cause analysis system has greatly reduced the amount of manual labor for operation staff and improves the user experience of its service. We also plan to deploy LEAST into more similar applications.

TABLE II
FLIGHT TICKET BOOKING EXAMPLE CASES IDENTIFIED BY LEAST.

Date	Identified Anomaly Path of Root Cause	Explainable Events
2019-11-19	Error in Step 3 ← Fare sources 3,9,16 ← Airline AC	Air Canada booking system unscheduled maintenance, which creates problems for a variety of different fare sources
2019-12-05	Error in Step 3 ← Airline SL ← Agent Office BKKOK275Q	Inaccurate data for Airline SL from Amadeus
2019-12-09	Error in Step 3 ← Error code P-1005-16-001-41203 Error in Step 3 ← Airline MU ← Fare source 5	Problem caused by internal system deployment
2020-01-23	Error in Step 1 ← Arrival city WUH	Lock-down of Wuhan City and many flights are cancelled
2020-02-15/20/28	Error in Step 1 ← Arrival city BKK	Australia extended travel ban from China and passengers sought Bangkok as a transfer point
2020-02-24	Error in Step 1 ← Departure city SEL Error in Step 1 ← Airline MU	COVID-19 broke out in South Korea

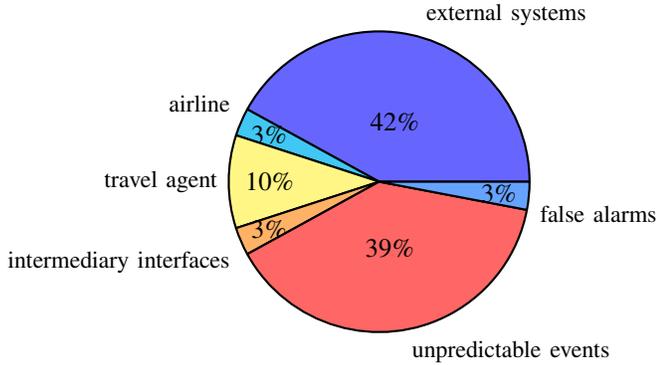


Fig. 7. Real-world evaluation results.

B. Gene Expression Data Analysis

We further apply our LEAST for gene expression data analysis, a traditional application of structure learning for BN [18], [26]. We use the three well-known datasets *Sachs* [29], *E. Coli* [27] and *Yeast* [27]. In each dataset, each node is a gene and each sample represents the normalized expression level values of some genes in a biological bench experiment. The underlying networks are commonly referred to as gene regulatory networks. We have the ground truth on the BN structure of them.

We run our LEAST method and the state-of-the-art NOTEARS method on a GPU. The evaluation results are reported in Table I. Regarding accuracy, we inspect the false direction rate (FDR), true positive rate (TPR), the false positive rate (FPR) and ROC-AUC metrics. We observe that the results of LEAST appear to be slightly better compared to NOTEARS on all gene datasets. LEAST detects 45 and 104 more true positive edges than NOTEARS on *E. Coli* and *Yeast*, respectively. Relatively, on *E. Coli*, LEAST improves the F_1 -score and AUC-ROC by 3.5% and 5.2%, respectively.

We conjecture that this is due to slight numerical instabilities due to the iterative nature of the process for computing $\bar{\delta}(\mathbf{W})$, which may eventually cause a higher variance in the final results. However, it appears that this additional source of noise acts as a form of implicit regularization, causing LEAST to exhibit higher accuracy in the gene dataset experiments.

On a GPU, in terms of the run time, the speed up ratio appears to be as not as significant as in the CPU case. This is could be due to the high level of parallelism possible in modern GPUs, concealing the speed up effects introduced by the lower complexity matrix operations in LEAST. However,

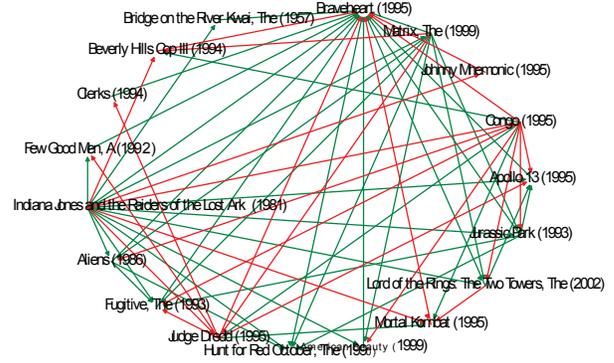


Fig. 8. Example extracted subgraph: Green and red edges indicate positive and negative learned weights, respectively.

note that, typical sizes of modern GPU memory does not allow for processing graphs larger than e.g. *Yeast*.

C. Recommendation Systems: MovieLens Case Study

We provide here a case study using the *MovieLens* dataset to show how structure learning might find an application in the context of recommendation systems. The purpose is to exemplify how structure learning, once it is possible to apply it on large-scale datasets, might play a role in understanding user behavior and building explainable recommendation systems. In particular, we find that the resulting DAG for this case study resembles item-to-item graphs typically encountered in neighborhood based recommendation systems [6]. Yet, since our method is agnostic to the real-world meaning of the random variables X_i , using structure learning one could also learn graphs which combine relations between items, users, item features and user features in a straight forward manner.

A qualitative inspection of the resulting DAG yields interesting insights into the data. We present some of the learned edges with the strongest weights in Table IV. We often find, that links with strong positive weights indicate very similar movies, i.e. movies from the same series, period or director.

Furthermore, a representative subgraph extracted around the Movie *Braveheart* is presented in Figure 8. One may directly interpret the learned structure as follows: given a user's rating for a specific movie, we start at the node corresponding to the movie i and follow outgoing edges to movies j , while multiplying the rating with weights for edges $i \rightarrow j$. If resulting values are positive one could predict whether based on the original rating for movie i , a user might like (positive

TABLE III
PROPERTIES OF REAL-WORLD DATASETS IN EXPERIMENTS.

Metric	Sachs		E. Coli		Yeast	
	NOTEARS	LEAST	NOTEARS	LEAST	NOTEARS	LEAST
# of Nodes	11		1,565		4,441	
# of Samples	1,000		1,565		4,441	
# of Exact Edges	17		3,648		12,873	
# of Predicted Edges	17	15	164	234	560	794
# of True Positive Edges	7	7	54	99	307	411
FDR	0.353	0.353	0.146	0.103	0.013	0.011
TPR	0.412	0.412	0.047	0.066	0.044	0.062
FPR	0.263	0.211	1.97×10^{-5}	1.97×10^{-5}	7.11×10^{-7}	9.14×10^{-7}
SHD	14	12	3,492	3,422	12,311	12,079
F_1 -score	0.412	0.437	0.073	0.108	0.082	0.119
AUC-ROC	0.925	0.947	0.831	0.883	0.891	0.919
GPU Time (Sec)	134	53	811	415	6,610	6,930

TABLE IV
EXAMPLES FOR TOP-10 LEARNED EDGES FOR *Movielens*

Link From		Link To	Weight	Remarks
Shrek 2 (2004)	→	Shrek (2001)	0.220	same series
Raiders of the Lost Ark (1981)	→	Star Wars: Episode IV (1977)	0.178	same main actor
Raiders of the Lost Ark (1981)	→	Indiana Jones and the Last Crusade (1989)	0.159	same series
Harry Potter and the Chamber of Secrets (2002)	→	Harry Potter and the Sorcerer’s Stone (2001)	0.159	same series
The Maltese Falcon (1941)	→	Casablanca (1942)	0.159	same period
Reservoir Dogs (1992)	→	Pulp Fiction (1994)	0.146	same director
North by Northwest (1959)	→	Rear Window (1954)	0.144	same director
Toy Story 2 (1999)	→	Toy Story (1995)	0.144	same series
Spider-Man (2002)	→	Spider-Man 2 (2004)	0.126	same series
Seven (1995)	→	The Silence of the Lambs (1991)	0.126	same genre

value) or dislike the movie j . This indicates the possibility of building a simple, yet explainable recommendation system based on the learned graph structure.

We also observe an interesting phenomenon related to the acyclicity constraint: So called “blockbuster” or very popular movies, which can be assumed to have been watched by the majority of users, such as e.g. *Star Wars: Episode V* (no outgoing, but 68 incoming links) or *Casablanca* (no outgoing, but 48 incoming links) tend to have a larger number of incoming links than outgoing links. On the other hand, less well known movies or movies indicative of a more specialized taste, such as *The New Land* (with no incoming, but 221 outgoing links), tend to have more outgoing links.

One possible explanation for this could be given along the thought: Since the majority of users enjoy e.g. *Star Wars*, we do not obtain any interesting insight into a user’s taste if we only knew that he or she would like e.g. *Star Wars: Episode V* (no outgoing, but 68 incoming links) or *Casablanca* (no outgoing, but 48 incoming links). On the other hand, enjoying a movie such as the 1972 Swedish movie *The New Land* (with no incoming, but 221 outgoing links) seems to be much more indicative of a user’s personal tastes in movies.

VII. CONCLUSIONS AND FUTURE WORK

Driven by a wide variety of business scenarios in Alibaba, we design LEAST, a highly accurate, efficient and scalable on structure learning algorithm for BN. LEAST is built upon a novel way to characterize graph acyclicity using an upper bound of the graph spectral radius. Benchmark evaluation exhibits that LEAST attains state-of-the-art performance in

terms of accuracy and efficiency. It is able to scale on BNs with hundreds of thousands nodes. Moreover, LEAST has been deployed in our production systems to serve a wide range of applications.

This work makes the first attempt to apply theoretically appealing method (BN) to address complex problems at the intersection of machine learning and causality, which are previously considered intractable. We believe that it would attract more efforts on applying structure learning for BN to problems involving large amount of random variables. We also hope to implement LEAST in a distributed environment to explore extremely large-scale applications.

REFERENCES

- [1] Remco R Bouckaert. Probabilistic network construction using the minimum description length principle. In *ECSQRU*, pages 41–48. Springer, 1993.
- [2] Kung-Ching Chang, Kelly Pearson, and Tan Zhang. Perron-frobenius theorem for nonnegative tensors. *CMS*, 6(2):507–520, 2008.
- [3] David Maxwell Chickering. Learning bayesian networks is np-complete. In *Learning from Data*, pages 121–130. Springer, 1996.
- [4] David Maxwell Chickering. Optimal structure identification with greedy search. *JMLR*, 3(Nov):507–554, 2002.
- [5] David Maxwell Chickering and David Heckerman. Efficient approximations for the marginal likelihood of bayesian networks with hidden variables. *ML*, 29(2-3):181–212, 1997.
- [6] Aaron Defazio and Tiberio Caetano. A graphical model formulation of collaborative filtering neighbourhood methods with fast maximum entropy training. *arXiv:1206.4622*, 2012.
- [7] Mathias Drton and Marloes H Maathuis. Structure learning in graphical modeling. *Annual Review of Statistics and Its Application*, 4:365–393, 2017.
- [8] Fei Fu and Qing Zhou. Learning sparse causal gaussian networks with experimental intervention: regularization and coordinate descent. *JASA*, 108(501):288–300, 2013.

- [9] Olivier Goudet, Diviyani Kalainathan, Philippe Caillou, Isabelle Guyon, David Lopez-Paz, and Michele Sebag. Learning functional causal models with generative neural networks. In *Explainable and Interpretable Models in CV and ML*, pages 39–80. Springer, 2018.
- [10] Jiaying Gu, Fei Fu, and Qing Zhou. Penalized estimation of directed acyclic graphs from discrete data. *Statistics and Computing*, 29(1):161–176, 2019.
- [11] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *TKDE*, 5(4):1–19, 2015.
- [12] David Heckerman, Dan Geiger, and David Chickering. Learning bayesian networks: The combination of knowledge and statistical data. *ML*, 20(3):197–243, 1995.
- [13] Alireza Sadeghi Hesar, Hamid Tabatabaee, and Mehrdad Jalali. Structure learning of bayesian networks using heuristic methods. In *ICIKM 2012*, 2012.
- [14] Patrik O Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In *NeurIPS*, pages 689–696, 2009.
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- [16] Jack Kuipers, Giusi Moffa, David Heckerman, et al. Addendum on the scoring of gaussian directed acyclic graphical models. *The Annals of Statistics*, 42(4):1689–1691, 2014.
- [17] Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural dag learning. *arXiv:1906.02226*, 2019.
- [18] Hao-Chih Lee, Matteo Danieletto, Riccardo Miotto, Sarah T Cherng, and Joel T Dudley. Scaling structural learning with no-bears to infer causal transcriptome networks. In *Pac Symp Biocomput*. World Scientific, 2019.
- [19] Szymon Łukasik and Sławomir Zak. Firefly algorithm for continuous constrained optimization tasks. In *ICCCI*, pages 97–106. Springer, 2009.
- [20] Young Woong Park and Diego Klabjan. Bayesian network learning via topological order. *JMLR*, 18(1):3451–3482, 2017.
- [21] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. MIT Press, 2017.
- [22] Joseph Ramsey, Madelyn Glymour, Ruben Sanchez-Romero, and Clark Glymour. A million variables and more: the fast greedy equivalence search algorithm for learning high-dimensional graphical causal models, with an application to functional magnetic resonance images. *IJDSA*, 3(2):121–129, 2017.
- [23] Robert W Robinson. Counting unlabeled acyclic digraphs. In *Com. Math. V*, pages 28–43. Springer, 1977.
- [24] Mauro Scanagatta, Giorgio Corani, Cassio P De Campos, and Marco Zaffalon. Learning treewidth-bounded bayesian networks with thousands of variables. In *NeurIPS*, pages 1462–1470, 2016.
- [25] Mauro Scanagatta, Cassio P de Campos, Giorgio Corani, and Marco Zaffalon. Learning bayesian networks with thousands of variables. In *NeurIPS*, pages 1864–1872, 2015.
- [26] Mauro Scanagatta, Antonio Salmerón, and Fabio Stella. A survey on bayesian network structure learning from data. *Progress in AI*, pages 1–15, 2019.
- [27] Thomas Schaffter, Daniel Marbach, and Dario Floreano. Genenetweaver: in silico benchmark generation and performance profiling of network inference methods. *Bioinformatics*, 27(16):2263–2270, 2011.
- [28] Mark Schmidt, Alexandru Niculescu-Mizil, Kevin Murphy, et al. Learning graphical model structure using l1-regularization paths. In *AAAI*, volume 7, pages 1278–1283, 2007.
- [29] Marco Scutari. Bn learning repository. In <https://www.bnlearn.com/bnrepository/>, 2020.
- [30] Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-gaussian acyclic model for causal discovery. *JMLR*, 7(Oct):2003–2030, 2006.
- [31] Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvärinen, Yoshinobu Kawahara, Takashi Washio, Patrik O Hoyer, and Kenneth Bollen. Directlingam: A direct method for learning a linear non-gaussian structural equation model. *JMLR*, 12(Apr):1225–1248, 2011.
- [32] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT Press, 2000.
- [33] Lambert M. Surhone, Mariam T. Tennoe, and Susan F. Henssonow. *Spectral Radius*. 2011.
- [34] Kostas Tzoumas, Amol Deshpande, and Christian S Jensen. Lightweight graphical models for selectivity estimation without independence assumptions. *PVLDB*, 4(11):852–863, 2011.
- [35] Xiangyu Wang, David Dunson, and Chenlei Leng. No penalty no tears: Least squares in high-dimensional linear models. In *ICML*, pages 1814–1822, 2016.
- [36] Sarah Wellen and David Danks. Learning causal structure through local prediction-error learning. In *CogSci*, volume 34, 2012.
- [37] Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks. In *ICML*, pages 7154–7163, 2019.
- [38] Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. In *NeurIPS*, pages 9472–9483, 2018.
- [39] Xun Zheng, Chen Dan, Bryon Aragam, Pradeep Ravikumar, and Eric P Xing. Learning sparse nonparametric dags. *arXiv:1909.13189*, 2019.
- [40] Shuheng Zhou. Thresholding procedures for high dimensional variable selection and statistical estimation. In *NeurIPS*, pages 2304–2312, 2009.
- [41] Shengyu Zhu and Zhitang Chen. Causal discovery with reinforcement learning. *arXiv:1906.04477*, 2019.

This figure "feizhu-graph.png" is available in "png" format from:

<http://arxiv.org/ps/2012.03540v1>

This figure "fig1.png" is available in "png" format from:

<http://arxiv.org/ps/2012.03540v1>

This figure "movie_lens_paper_1.png" is available in "png" format from:

<http://arxiv.org/ps/2012.03540v1>