

Fast Bayesian Network Structure Learning using Quasi-Determinism Screening

Thibaud Rahier, Sylvain Marié, Stéphane Girard, Florence Forbes

► **To cite this version:**

Thibaud Rahier, Sylvain Marié, Stéphane Girard, Florence Forbes. Fast Bayesian Network Structure Learning using Quasi-Determinism Screening. 2018. <hal-01691217v3>

HAL Id: hal-01691217

<https://hal.archives-ouvertes.fr/hal-01691217v3>

Submitted on 12 Mar 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Fast Bayesian Network Structure Learning using Quasi-Determinism Screening

Thibaud Rahier
Univ. Grenoble Alpes, Inria
& Schneider Electric
Grenoble, France

Sylvain Marie
Schneider Electric
Grenoble, France

Stephane Girard
Univ. Grenoble Alpes, Inria
Grenoble, France

Florence Forbes
Univ. Grenoble Alpes, Inria
Grenoble, France

Abstract

Learning the structure of Bayesian networks from data is a NP-Hard problem that involves optimization over a super-exponential sized space. In this work, we show that in most real life datasets, a number of the arcs contained in the final structure can be pre-screened at low computational cost with a limited impact on the global graph score. We formalize the identification of these arcs via the notion of quasi-determinism, and propose an associated algorithm that narrows the structure learning task down to a subset of the original variables. We show, on diverse benchmark datasets, that this algorithm exhibits a significant decrease in computational time and complexity for only a little decrease in performance score.

1 INTRODUCTION

Bayesian networks are probabilistic graphical models that present interest both in terms of knowledge discovery and density estimation. Learning Bayesian networks from data has been however proven to be NP-Hard by Chickering (1996).

There has been extensive work on tackling the ambitious problem of Bayesian network structure learning from observational data. Algorithms fall under two main categories: *constraint-based* and *score-based*.

Constraint-based structure learning algorithms rely on testing for conditional independencies that hold in the data in order to reconstruct a Bayesian network encoding these independencies. The *PC* algorithm by Spirtes et al. (2000) was the first practical application of this idea, followed by several optimized approaches as the *fast incremental association* (Fast-IAMB) algorithm from Yaramakala and Margaritis (2005).

Score-based structure learning relies on the definition of a network score, then on the search for the best-scoring structure among all possible directed acyclic graphs (DAGs). The number of possible DAG structures with n nodes is of order $2^{\mathcal{O}(n^2)}$, which prevents exhaustive search when n is typically larger than 30, even for the most recent algorithms Silander and Myllymaki (2012) or Yuan et al. (2013).

Most of the score-based algorithms used in practice therefore rely on heuristics, as the original approach from Cooper and Herskovits (1992) which supposed a prior ordering of the variables to perform parent set selection, or Bouckaert (1995) who proposed to search through the structure space using greedy hill climbing with random restarts. Since these first algorithms, different approaches have been proposed: some based on the search for an optimal ordering as Chen et al. (2008) or Teyssier and Koller (2012), others on optimizing the search task in accordance to a given score (BIC or BDe) as de Campos and Ji (2011) or Scanagatta et al. (2015).

Meanwhile, data itself may contain determinism, for example in the fields of cancer risk identification (de Morais et al. (2008)) or nuclear safety (Mabrouk et al. (2014)). Moreover, data is increasingly collected and generated by software systems whether in social networks, smart buildings, smart grid, smart cities or the internet of things (IoT) in general (Koo et al. (2016)). These systems in their vast majority rely on relational data models or lately on semantic data models (El Kaed et al. (2016)) which cause deterministic relationships between variables to be more and more common in datasets. Determinism has been shown to interfere with Bayesian network structure learning, notably constraint-based methods as mentioned by Luo (2006).

In this paper, we focus on score-based algorithms. After reminding the background of Bayesian network structure learning in section 2, we state some theoretical results in section 3, that enable to bridge the gap between deter-

minism and Bayesian network scoring.

In section 4, exploiting the intuition brought by these theoretical results, we propose and study the complexity of the *quasi deterministic screening* algorithm. The idea is that some of the arcs contained in the desired Bayesian network can be learned during a quick screening phase where quasi-deterministic relationships are detected, thus narrowing the learning task down to a subset of the original variables.

In practice, not only does this algorithm accelerate the overall learning procedure with very low performance loss, but it also leads to sparser and therefore more interpretable graphs than state of the art methods, as presented using benchmark datasets in section 5.

Finally, section 6 is dedicated to a discussion and to numerous perspectives emerging from this work. Proofs of all lemmas and propositions are available in the supplementary material.

2 BAYESIAN NETWORK STRUCTURE LEARNING

2.1 Bayesian networks

Let $\mathbf{X} = (X_1, \dots, X_n)$ be a n -tuple of categorical random variables with respective value sets $Val(X_1), \dots, Val(X_n)$. The distribution of \mathbf{X} is denoted by, $\forall \mathbf{x} = (x_1, \dots, x_n) \in Val(\mathbf{X})$,

$$p(\mathbf{x}) = P(X_1 = x_1, \dots, X_n = x_n).$$

For $I \subset \llbracket 1, n \rrbracket$, we define $\mathbf{X}_I = \{X_i\}_{i \in I}$, and the notation $p(\cdot)$ and $p(\cdot|\cdot)$ is extended to the marginals and conditionals of any subset of variables: $\forall(\mathbf{x}_I, \mathbf{x}_J) \in Val(\mathbf{X}_{I \cup J})$, $p(\mathbf{x}_I|\mathbf{x}_J) = P(\mathbf{X}_I = \mathbf{x}_I | \mathbf{X}_J = \mathbf{x}_J)$.

Moreover, we suppose that D is a dataset containing M *i.i.d.* instances of (X_1, \dots, X_n) . All quantities empirically computed from D will be written with a \cdot^D exponent (e.g. p^D refers to the empirical distribution with respect to D). Finally, D_I refers to the restriction of D to the observations of \mathbf{X}_I .

A Bayesian network is an object $\mathcal{B} = (G, \theta)$ where

- $G = (V, A)$ is a directed acyclic graph (DAG) structure with V the set of nodes and $A \subset V \times V$ the set of arcs. We suppose $V = \llbracket 1, n \rrbracket$ where each node $i \in V$ is associated with the random variable X_i , and $\pi^G(i) = \{j \in V \text{ s.t. } (j, i) \in A\}$ is the set of i 's parents in G . The exponent G may be dropped for clarity when the referred graph is obvious.
- $\theta = \{\theta_i\}_{i \in V}$ is a set of parameters. Each θ_i defines the local conditional distribution $P(X_i | \mathbf{X}_{\pi(i)})$. More precisely, $\theta_i = \{\theta_{x_i | \mathbf{x}_{\pi(i)}}\}$ where for $i \in$

$$V, x_i \in Val(X_i) \text{ and } \mathbf{x}_{\pi(i)} \in Val(\mathbf{X}_{\pi(i)}),$$

$$\theta_{x_i | \mathbf{x}_{\pi(i)}} = p(x_i | \mathbf{x}_{\pi(i)}).$$

A Bayesian network $\mathcal{B} = (G, \theta)$ encodes the following factorization of the distribution of \mathbf{X} : for $\mathbf{x} = (x_1, \dots, x_n) \in Val(\mathbf{X})$,

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | \mathbf{x}_{\pi^G(i)}) = \prod_{i=1}^n \theta_{x_i | \mathbf{x}_{\pi^G(i)}}.$$

Such a factorization notably implies that *each variable is independent of its non-descendants given its parents*.

2.2 Score-based approach to Bayesian network structure learning

Suppose we have a *scoring* function $s : DAG_V \rightarrow \mathbb{R}$, where DAG_V is the set of all possible DAG structures with node set V . Score-based Bayesian network structure learning comes down to solving the following combinatorial optimization problem:

$$G^* \in \operatorname{argmax}_{G \in DAG_V} s(G). \quad (1)$$

It can be shown that $2^{\frac{n(n-1)}{2}} \leq |DAG_V| \leq 2^{n(n-1)}$ where $|V| = n$. There are therefore $2^{\mathcal{O}(n^2)}$ possible DAG structures containing n nodes: the size of DAG_V is said to be super-exponential in $|V|$.

Most scoring functions used in practice are based on the likelihood function. The most straightforward being the Max log-likelihood score, that we now present.

The Max log-likelihood score Let $l(\theta : D) = \log(p_\theta(D))$ be the log-likelihood of the set of parameters θ given the dataset D .

For a given DAG structure $G \in DAG_V$, we define the Max log-likelihood (MLL) score of G associated with a dataset D as:

$$s^{MLL}(G : D) = \max_{\theta \in \Theta_G} l(\theta : D).$$

where Θ_G is the set of all θ 's such that $\mathcal{B} = (G, \theta)$ is a Bayesian network.

The MLL score is very straightforward and intuitive, but it favors denser structures: if $G_1 = (V, A_1)$ and $G_2 = (V, A_2)$ are two graph structures such that $A_1 \subset A_2$, we can show that: $s^{MLL}(G_1 : D) \leq s^{MLL}(G_2 : D)$.

There are two main (non-exclusive) approaches to solve this problem:

- Constrain the structure space to avoid learning overly complex graphs, which is the idea of hybrid structure learning algorithms such as the sparse candidate algorithm presented by Friedman et al. (1999), or the Max-Min Hill Climbing (MMHC) algorithm introduced by Tsamardinos et al. (2006).

- Use a score that induces a goodness-of-fit versus complexity tradeoff, such as:
 - BIC (Schwarz et al. (1978)), which is a penalized version of the MLL score.
 - BDe (Heckerman et al. (1995)), which is derived from the marginalization of the likelihood, implying a prior implicitly penalizing the model’s complexity.

The BDe score In this paper, we will use the BDe score to evaluate a BN structure’s quality, as it is done in several recent papers as Teyssier and Koller (2012), or Nie et al. (2016). This score is known to be a good indicator of generalization performance.

The BDe score of $G \in DAG_V$ is defined as the log of the marginal likelihood, integrated against a Dirichlet prior. Assuming a uniform prior over all network structures, we have

$$s^{BDe}(G : D) = \log \left(\int_{\theta \in \Theta_G} \underbrace{p(D|\theta, G)}_{\text{Likelihood}} \underbrace{p(\theta|G)}_{\text{Dirichlet prior}} d\theta \right)$$

In practice, we often use the BDeu score, introduced by Buntine (1991), which is a particular case of BDe where the Dirichlet prior assigns the same probability to all configurations of $X_i \cup X_{\pi(i)}$ for each i .

3 DETERMINISM AND BAYESIAN NETWORKS

3.1 Definitions

We propose the following definitions of determinism and deterministic DAGs using the notion of conditional entropy. In this paper, determinism will always be meant empirically, with respect to a dataset D .

Definition 1 Determinism wrt D

Given a dataset D containing observations of variables X_i and X_j , the relationship $X_i \rightarrow X_j$ is deterministic with respect to D iff $H^D(X_i|X_j) = 0$, where

$$H^D(X_i|X_j) = - \sum_{x_i, x_j} p^D(x_i, x_j) \log(p^D(x_i|x_j))$$

is the empirical conditional Shannon entropy.

It is straightforward to prove that Definition 1 relates to a common and intuitive perception of determinism, as the one presented by Luo (2006). Indeed,

$$H^D(X_i|X_j) = 0 \Leftrightarrow \forall x_j \in Val(X_j), \exists! x_i \in Val(X_i) \text{ s.t. } p^D(x_i|x_j) = 1.$$

This definition is naturally extended to \mathbf{X}_I and \mathbf{X}_J for $I, J \subset V$, i.e. $\mathbf{X}_I \rightarrow \mathbf{X}_J$ is deterministic with respect to D iff $H^D(\mathbf{X}_J|\mathbf{X}_I) = 0$.

Definition 2 Deterministic DAG wrt D

$G \in DAG_V$ is said to be deterministic with respect to D iff $\forall i \in V$ s.t. $\pi^G(i) \neq \emptyset$, $\mathbf{X}_{\pi^G(i)} \rightarrow X_i$ is deterministic wrt D .

3.2 Deterministic trees and MLL score

We first recall a lemma that relates the MLL score presented in section 2 to the notion of empirical conditional entropy. This result is well known and notably stated by Koller and Friedman (2009).

Lemma 1 For $G \in DAG_V$ associated with variables X_1, \dots, X_n observed in a dataset D ,

$$s^{MLL}(G : D) = -M \sum_{i=1}^n H^D(X_i|\mathbf{X}_{\pi(i)})$$

where by convention $H^D(X_i|\mathbf{X}_\emptyset) = H^D(X_i)$.

The next proposition follows then straightforwardly. We remind that a tree is a DAG in which each node has exactly one parent, except for the root node which has none.

Proposition 1 If T is a deterministic tree with respect to D , then T is a solution of (1):

$$s^{MLL}(T : D) = \max_{G \in DAG_V} s^{MLL}(G : D).$$

It is worth noticing that complete DAGs also maximize the MLL score. The main interest of Proposition 1 resides in the fact that, under the (strong) assumption that a deterministic tree exists, we are able to explicit a sparse solution of (1), with $n - 1$ arcs instead of $\frac{n(n-1)}{2}$ for a complete DAG.

3.3 Deterministic forests and the MLL score

The deterministic tree assumption of Proposition 1 is very restrictive. In this section, it is extended to deterministic forests, defined as follows:

Definition 3 Deterministic forest wrt D

$F \in DAG_V$ is said to be a deterministic forest with respect to D iff $F = \bigcup_{k=1}^p T_k$, where T_1, \dots, T_p are p disjoint deterministic trees wrt $D_{V_{T_1}}, \dots, D_{V_{T_p}}$ respectively and s.t. $\bigcup_{k=1}^p V_{T_k} = V$.

In the expression $\bigcup_{k=1}^p T_k$, \cup is the canonical union for graphs: $G \cup G' = (V_G \cup V_{G'}, A_G \cup A_{G'})$.

For a given deterministic forest F wrt D , we define

$$R(F) = \{i \in V \mid \pi^F(i) = \emptyset\}$$

the set of F 's roots (the union of the roots of each of its trees).

Proposition 2 *Suppose F is a deterministic forest wrt D . Let $G_{R(F)}^*$ be a solution of the structure learning optimization problem (1) for $\mathbf{X}_{R(F)}$ and the MLL score i.e.*

$$s^{MLL}(G_{R(F)}^* : D_{R(F)}) = \max_{G \in DAG_{R(F)}} s^{MLL}(G : D_{R(F)}).$$

Then, $G^* = F \cup G_{R(F)}^*$ is a solution of (1) for \mathbf{X} , i.e.

$$s^{MLL}(G^* : D) = \max_{G \in DAG_V} s^{MLL}(G : D).$$

As opposed to Proposition 1, the assumptions of Proposition 2 are always formally verified: if there is no determinism in the dataset D , then $R(F) = V$, and every tree T_k is formed of a single root node. In that case, solving problem (1) for $G_{R(F)}^*$ is the same as solving it for G^* . Of course, we are interested in the case where $R(F) < n$, as this enables us to focus on a smaller structure learning problem while still having the guarantee to learn the optimal Bayesian network with regards to the MLL score.

As seen in section 2, the main issue with the MLL score is that it favors complete graphs. However, a deterministic forest F containing p trees is very sparse ($n - p$ arcs), and even if the graph $G_{R(F)}^*$ is dense, the graph $G^* = F \cup G_{R(F)}^*$ may still satisfy sparsity conditions.

4 STRUCTURE LEARNING WITH QUASI-DETERMINISM SCREENING

4.1 Quasi-determinism

When it comes to Bayesian network structure learning algorithms, even heuristics are computationally intensive. We would like to use the theoretical results presented in section 3 to simplify the structure learning problem.

Our idea is to narrow the structure learning problem down to a subset of the original variables: the roots of a deterministic forest, in order to significantly decrease the overall computation time. This is what we call determinism screening.

However, in many datasets, one does not observe real empirical determinism, although there are very strong

relationships between some of the variables. We therefore propose to *relax* the aforementioned determinism screening to quasi-determinism screening, where *quasi* is meant with respect to a parameter ϵ : we talk about ϵ -quasi-determinism.

There are several ways to measure how close a relationship is from deterministic. Huhtala et al. (1999) consider the minimum number of observations that must be dropped from the data for the relationship to be deterministic. Since we are in a score-maximization context, we will rather use ϵ as a threshold on the empirical conditional entropy. The following definition is the natural generalization of Definition 1.

Definition 4 ϵ -quasi-determinism (ϵ -qd)

Given a dataset D containing observations of variables X_i and X_j , the relationship $X_i \rightarrow X_j$ is ϵ -qd wrt D iff $H^D(X_j|X_i) \leq \epsilon$.

It has been seen in Proposition 2 that a deterministic forest is the subgraph of an optimal DAG with respect to the MLL score, while still satisfying sparsity conditions. Such a forest is therefore very promising with regards to the fit-complexity tradeoff (typically evaluated by scores such as BDe or BIC).

Combining this intuition with the ϵ -qd criteria presented in Definition 4, we propose the quasi-determinism screening approach to Bayesian network structure learning, defined in the next subsections.

4.2 Quasi-determinism screening algorithm

Algorithm 1 details how to find the simplest ϵ -qd forest F_ϵ from a dataset D and a threshold ϵ . Here *simplest* refers to the complexity in terms of number of parameters.

This algorithm takes for input:

- D : a dataset containing M observations of \mathbf{X} ,
- ϵ : a threshold for quasi-determinism.

The next proposition states that Algorithm 1 is well defined.

Proposition 3 *For any rightful input D and ϵ , the output of Algorithm 1 is a forest (i.e. a directed acyclic graph with at most one parent per node).*

4.3 Learning Bayesian networks using quasi-determinism screening

We now present Algorithm 2, which uses quasi-determinism screening to accelerate Bayesian network structure learning. This algorithm takes the following input:

Algorithm 1 Quasi-determinism screening (qds)

Input: D, ϵ

- 1: Compute empirical conditional entropy matrix $\mathbb{H}^D = (H^D(X_i|X_j))_{1 \leq i, j \leq n}$
- 2: **for** $i = 1$ to n **do** #identify the set of potential ϵ -qd parents for each i
- 3: compute $\pi_\epsilon(i) = \{j \in \llbracket 1, n \rrbracket \setminus \{i\} \mid \mathbb{H}_{ij}^D \leq \epsilon\}$
- 4: **for** $i = 1$ to n **do** #check for cycles in ϵ -qd relations
- 5: **if** $\exists j \in \pi_\epsilon(i)$ s.t. $i \in \pi_\epsilon(j)$ **then**
- 6: **if** $\mathbb{H}_{ij}^D \leq \mathbb{H}_{ji}^D$ **then**
- 7: $\pi_\epsilon(j) \leftarrow \pi_\epsilon(j) \setminus \{i\}$
- 8: **else**
- 9: $\pi_\epsilon(i) \leftarrow \pi_\epsilon(i) \setminus \{j\}$
- 10: **for** $i = 1$ to n **do** #choose the simplest among all potential parents
- 11: $\pi_\epsilon^*(i) \leftarrow \underset{j \in \pi_\epsilon(i)}{\operatorname{argmin}} |Val(X_j)|$
- 12: Compute forest $F_\epsilon = (V_{F_\epsilon}, A_{F_\epsilon})$ where $V_{F_\epsilon} = \llbracket 1, n \rrbracket$ and $A_{F_\epsilon} = \{(\pi_\epsilon^*(i), i) \mid i \in \llbracket 1, n \rrbracket \text{ s.t. } \pi_\epsilon^*(i) \neq \emptyset\}$

Output: F_ϵ

- D : a dataset containing M observations of \mathbf{X} ,
- ϵ : a threshold for quasi-determinism,
- *sota-BNSL*: a state of the art structure learning algorithm, taking for input a dataset, and returning a Bayesian network structure.

Algorithm 2 Bayesian network structure learning with quasi deterministic screening (qds-BNSL)

Input: $D, \epsilon, \text{sota-BNSL}$

- 1: Compute F_ϵ by running **Algorithm 1** with input D and ϵ
- 2: Identify $R(F_\epsilon) = \{i \in \llbracket 1, n \rrbracket \mid \pi^{F_\epsilon}(i) = \emptyset\}$, the set of F_ϵ 's roots.
- 3: Compute $G_{R(F_\epsilon)}^*$ by running *sota-BNSL* on $D_{R(F_\epsilon)}$
- 4: $G_\epsilon^* \leftarrow F_\epsilon \cup G_{R(F_\epsilon)}^*$

Output: G_ϵ^*

The extension of the definition of determinism to quasi-determinism (Definition 3) prevents us to have ‘hard’ guarantees as those presented in Proposition 2. However, we are able to explicit bounds for the MLL score of a graph G_ϵ^* returned by Algorithm 2, as stated in the following Proposition.

Proposition 4 *Let ϵ, D and *sota-BNSL* be rightful input to Algorithm 2, and G_ϵ^* the associated output. Then, if *sota-BNSL* is exact (i.e. always returns an optimal solution) with respect to the MLL score, we have*

the following lower bound for $s^{MLL}(G_\epsilon^* : D)$:

$$s^{MLL}(G_\epsilon^* : D) \geq \left(\max_{G \in DAG_V} s^{MLL}(G : D) \right) - Mn\epsilon.$$

In practice, this bound is not very tight and this result therefore has few applicative potential. However, it shows that:

$$s^{MLL}(G_\epsilon^* : D) \xrightarrow{\epsilon \rightarrow 0} \max_{G \in DAG_V} s^{MLL}(G : D).$$

In other words, $\epsilon \mapsto s^{MLL}(G_\epsilon^* : D)$ is continuous in 0, and Proposition 4 generalizes Proposition 2.

Algorithm 2 is promising, notably if for small ϵ we can significantly decrease the number of variables to be considered by *sota-BNSL*. We would in this case have a much faster algorithm with a controlled performance loss.

4.4 Complexity analysis

Complexity of the state of the art algorithm The number of possible DAG structures being super exponential in the number of nodes, state of the art algorithms do not entirely explore the structure space but use smart caching and pruning methods to have a good performance & computation time trade-off.

Let *sota-BNSL* be a state of the art Bayesian network structure learning algorithm and $C_{\text{sota}}(M, n)$ be its complexity.

$C_{\text{sota}}(M, n)$ should typically be thought of as linear in M and exponential, or at least high degree polynomial, in n for the best algorithms.

Complexity of Algorithm 1 We have the following decomposition of the complexity of Algorithm 1:

1. Lines 1-3: $\mathcal{O}(Mn^2)$. Computation of \mathbb{H}^D : we need counts for every couple (X_i, X_j) for $i < j$ (each time going through all rows of D), which implies $M \frac{n(n-1)}{2}$ operations.
2. lines 4-9: $\mathcal{O}(n^2)$. Going through all elements of \mathbb{H}^D once.
3. lines 10-12: $\mathcal{O}(n^2)$. Going through all elements of \mathbb{H}^D once.

Overall one has that $C_{\text{Alg1}}(M, n) = \mathcal{O}(Mn^2)$.

Complexity of Algorithm 2 For a given dataset D , we define:

$$\forall \epsilon \geq 0, n_r(\epsilon) = |R(F_\epsilon)|.$$

The function $n_r(\cdot)$, associates to $\epsilon \geq 0$ the number of roots of the forest F_ϵ returned by Algorithm 1. The complexity of Algorithm 2 then decomposes as:

1. Line 1: $\mathcal{O}(Mn^2)$. Run of Algorithm 1.
2. Lines 2-4: $C_{sota}(M, n_r(\epsilon))$. Run of *sota-BNSL* on reduced dataset $D_{R(F_\epsilon)}$ with $n_r(\epsilon)$ columns.

This yields $C_{Alg2}(M, n) = \mathcal{O}(Mn^2) + C_{sota}(M, n_r(\epsilon))$.

We are interested in how much it differs from $C_{sota}(M, n)$, which depends mainly on:

- how $n_r(\epsilon)$ compares to n ,
- how $C_{sota}(M, n)$ varies with respect to n .

It is hard to obtain a theoretical difference $C_{sota} - C_{Alg2}$, since it is not clear how to estimate the complexity of state of the art learning algorithms, often based on local search heuristics. However, we know that all Bayesian network structure learning algorithms are very time-intensive, which lets us expect an important decrease in computational time for Algorithm 2 compared to a state of the art algorithm.

In the next section, we run a state of the art structure learning algorithm and Algorithm 2 on benchmark datasets in order to confirm this intuition.

5 EXPERIMENTS

5.1 Experimental setup

Data Table 1 summarizes the data used in our experiments. We considered the largest open-source categorical datasets among those presented¹ by Davis and Domingos (2010) and available on the UCI repository (Dheeru and Karra Taniskidou (2017)): 20 Newsgroup, Adult, Book, Covertypes, KDDCup 2000, MSNBC, MSWeb, Plants, Reuters-52 and USCensus. Moreover, as it was done by Scanagatta et al. (2016), we chose the largest Bayesian networks available in the literature², for each of which we simulated 10000 observations: Andes, Hailfinder, Hepar 2, Link, Munin 1-4, PathFinder and Win95pts.

Programmation details and choice of *sota-BNSL*

Most of the code associated with this project was done in R, enabling an optimal exploitation of the `bnlearn` package from Scutari (2009), which is a very good reference among open-source packages dealing with Bayesian networks structure learning.

We need a state of the art Bayesian network structure learning algorithm to obtain a baseline performance. After carefully evaluating several algorithms implemented in the `bnlearn` package, we chose to use *Greedy Hill*

¹<http://alchemy.cs.washington.edu/papers/davis10a/>

²<http://www.bnlearn.com/bnrepository/>

Table 1: Datasets presentation

| name | short name | n | M |
|----------------|------------|------|---------|
| 20 newsgroups | 20ng | 930 | 11293 |
| adult | adult | 125 | 36631 |
| book | book | 500 | 8700 |
| covertime | covertime | 84 | 30000 |
| kddcup 2000 | kddcup | 64 | 180092 |
| msnbc | msnbc | 17 | 291326 |
| msweb | msweb | 294 | 29441 |
| plants | plants | 69 | 17412 |
| reuters 52 | r52 | 941 | 6532 |
| uscensus | uscensus | 68 | 2458285 |
| andes | andes | 223 | 10000 |
| hailfinder | hailfinder | 56 | 10000 |
| hepar 2 | hepar2 | 70 | 10000 |
| link | link | 724 | 10000 |
| munin 1 | munin1 | 186 | 10000 |
| munin 2 | munin2 | 1003 | 10000 |
| munin 3 | munin3 | 1041 | 10000 |
| munin 4 | munin4 | 1038 | 10000 |
| pathfinder | pathfinder | 109 | 10000 |
| windows 95 pts | win95pts | 76 | 10000 |

Climbing with random restarts and a tabu list, as it consistently outperformed other built-in algorithms both in time and score, in addition to being also used as a benchmark algorithm in the literature, notably by Teyssier and Koller (2012). In this section, we refer to this algorithm as *sota-BNSL*.

Choice of ϵ for *qds-BNSL* An approach to choosing ϵ in the case of the *qds-BNSL* algorithm is to pick values for $n_r(\epsilon)$, and manually find the corresponding values for ϵ . For a given dataset and $x \in [0, 1]$, we define $\epsilon_x = n_r^{-1}(\lfloor xn \rfloor)$. In other words, ϵ_x is the value of ϵ for which the number of roots of the qd forest F_ϵ represents a proportion x of the total number of variables.

The computation of ϵ_x is not problematic: once \mathbb{H}^D is computed and stored, evaluating $n_r(\epsilon)$ is done in constant time, and finding one of $n_r(\cdot)$'s quantiles is doable in $\mathcal{O}(\log(n))$ operations (dichotomy), which is negligible compared to the overall complexity of the screening.

Algorithm evaluation The algorithms are evaluated using 3 axes of performance:

- Bayesian network BDeu score presented in section 2 with equivalent sample size (ESS) equal to 5, inspired from Teyssier and Koller (2012).
- Number of arcs of the learned Bayesian network. The BDeu score naturally penalizes overly complex models (in terms of number of parameters), it is

however interesting to look at the number of arcs, as it is a straightforward way to evaluate how complex a Bayesian network structure appears to a human expert (and therefore how interpretable this structure is).

- Computing time t_{run} (all algorithms were run on the same machine).

It is essential to remark that *sota-BNSL* is used both to obtain a baseline performance and inside *qds-BNSL*. In both cases, it is run with the same settings until convergence. The comparison of computing times is therefore fair.

We present the obtained results for our selected state of the art algorithm *sota-BNSL*, and 3 versions of *qds-BNSL*. For each dataset, we selected $\epsilon \in \{\epsilon_{0.9}, \epsilon_{0.75}, \epsilon_{0.5}\}$, corresponding to a restriction of *sota-BNSL* to 90%, 75% and 50% of the original variables respectively.

The results are shown in Tables 2-4, one per evaluation criterion. In each table, the actual value of the criterion is displayed for *sota-BNSL* (*sota*), and the relative difference is displayed for the three versions of *qds-BNSL* we consider ($\mathbf{qds}_{\epsilon_{0.9}}$, $\mathbf{qds}_{\epsilon_{0.75}}$ and $\mathbf{qds}_{\epsilon_{0.5}}$).

5.2 Results

Score It appears in Table 2 that the decrease in BDeu score is smaller than 5% for all the considered datasets when 90% of the variables remain after the pre-screening ($\mathbf{qds}_{\epsilon_{0.9}}$), and for most of them when 75% of the variables remain ($\mathbf{qds}_{\epsilon_{0.75}}$). This is also observed with $\epsilon_{0.5}$ for datasets that contain a lot of very strong pairwise relationships as *kddcup*, *msweb*, or *munin2-4*.

Computing time Table 3 shows a significant decrease in computational time for *qds-BNSL*, which is all the more important that ϵ is large. In the best cases, we have both a very small decrease in BDeu score, and an important decrease in computational time. For example, the algorithm *qds-BNSL* with $\epsilon = \epsilon_{0.5}$ is 55% faster for *msweb*, and 54% for *munin 3*, while implying only around 1% decrease in score compared to *sota-BNSL*. If we allow a 5% score decrease, *qds-BNSL* can be up to 70% faster (20 newsgroups, book, *msnbc*, *kddcup*, *hepar2*, *pathfinder*).

These results confirm the complexity analysis of the previous section, in which we supposed that the screening phase had a very small computational cost compared to the standard structure learning phase.

Complexity As showed by Table 4, Bayesian networks learnt with *qds-BNSL* are consistently less complex than those learnt with *sota-BNSL*. Several graphs learnt with

Table 2: BDeu score averaged by observation. Every result that is less than 5% smaller than *sota-BNSL*’s score is boldfaced.

| dataset | sota | $\mathbf{qds}_{\epsilon_{0.9}}$ (%) | $\mathbf{qds}_{\epsilon_{0.75}}$ (%) | $\mathbf{qds}_{\epsilon_{0.5}}$ (%) |
|------------|---------|--|---|--|
| 20ng | -142.71 | -0.66 | -2.13 | -4.78 |
| adult | -12.86 | -0.16 | -0.05 | -4.01 |
| book | -34.81 | -0.80 | -1.69 | -4.64 |
| covertime | -13.60 | -0.21 | -1.23 | -11.7 |
| kddcup | -2.38 | -0.31 | -1.04 | -3.83 |
| msnbc | -6.19 | -0.14 | -2.62 | -4.64 |
| msweb | -9.77 | +0.03 | -0.07 | -0.99 |
| plants | -13.03 | -2.57 | -7.56 | -20.92 |
| r52 | -95.48 | -0.76 | -1.96 | -6.11 |
| uscensus | -23.20 | -0.27 | -1.75 | -10.39 |
| andes | -93.23 | -0.49 | -6.22 | -16.57 |
| hailfinder | -49.63 | -0.06 | -2.71 | -10.21 |
| hepar2 | -32.60 | -0.28 | -1.36 | -3.22 |
| link | -215.68 | +0.10 | +1.10 | -16.99 |
| munin1 | -41.15 | -0.09 | -0.16 | -9.95 |
| munin2 | -171.82 | -0.02 | -0.02 | -1.83 |
| munin3 | -165.09 | 0.00 | 0.00 | -1.10 |
| munin4 | -186.11 | -0.02 | -0.02 | -3.86 |
| pathfinder | -26.65 | -0.66 | -0.70 | -4.88 |
| win95pts | -9.22 | +0.06 | -1.08 | -9.15 |

$\mathbf{qds}_{\epsilon_{0.5}}$ are more than 30% sparser while still scoring less than 5% below state of the art: 20 newsgroups, book, *kddcup 2000*, *msnbc*, *msweb* and *hepar 2*.

Figure 1 and 2 display two Bayesian networks learnt on the ‘*msnbc*’ dataset. They provide an interesting example of the sparsity induced by *qds-BNSL*. After the $\mathbf{qds}_{\epsilon_{0.5}}$ -screening phase, half of the variables (corresponding to the nodes in white) are considered to be sufficiently explained by V_1 . They are therefore not taken into account by *sota-BNSL*, which is run only on the variables corresponding to the nodes in gray.

In the case of *msnbc*, this restriction of the learning problem implies only a small decrease in the final graph’s generalization performance (as reflected by the BDeu scores), while being 7 times faster to compute and enabling a significantly better readability.

In this processed version of the *msnbc* dataset (Davis and Domingos (2010)), each variable contains a binary information regarding the visit of a given page from the *msnbc.com* website³. The Bayesian network displayed in Figure 2 shows in a compact way the influence be-

³more details: <http://archive.ics.uci.edu/ml/machine-learning-databases/msnbc-ml/msnbc.data.html>

Table 3: Computation time (seconds). Every result that corresponds to a BDeu score less than 5% smaller than *sota-BNSL*’s score is boldfaced.

| dataset | sota (seconds) | qds $_{\epsilon_{0.9}}$ (%) | qds $_{\epsilon_{0.75}}$ (%) | qds $_{\epsilon_{0.5}}$ (%) |
|------------|-------------------|--------------------------------|---------------------------------|--------------------------------|
| 20ng | 21,495 | -1.62 | -42.66 | -72.94 |
| adult | 1,02 | -6.61 | -22.03 | -61.20 |
| book | 7,600 | -23.61 | -40.33 | -71.30 |
| covertypes | 565 | -6.80 | -33.22 | -71.13 |
| kddcup | 2,167 | -11.49 | -32.85 | -73.59 |
| msnbc | 252 | -20.66 | -60.61 | -85.65 |
| msweb | 4,701 | -6.29 | -9.86 | -55.08 |
| plants | 455 | -46.93 | -61.93 | -84.07 |
| r52 | 18,630 | -13.58 | -38.47 | -76.71 |
| uscensus | 21,782 | -0.44 | -31.54 | -77.68 |
| andes | 898 | -2.23 | -27.42 | -69.91 |
| hailfinder | 46 | -5.31 | -17.46 | -54.71 |
| hepar2 | 76 | -4.05 | -42.56 | -70.00 |
| link | 7,240 | -12.03 | -10.58 | -61.30 |
| munin1 | 497 | -7.42 | -17.23 | -59.14 |
| munin2 | 7,093 | -20.46 | -21.66 | -43.68 |
| munin3 | 11,558 | -36.91 | -29.20 | -54.19 |
| munin4 | 8,550 | -7.87 | -13.08 | -39.06 |
| pathfinder | 231 | -14.01 | -35.38 | -69.48 |
| win95pts | 132 | -6.05 | -31.41 | -69.07 |

Table 4: Networks’ number of arcs. Every result that corresponds to a BDeu score less than 5% smaller than *sota-BNSL*’s score is boldfaced.

| dataset | sota | qds $_{\epsilon_{0.9}}$ (%) | qds $_{\epsilon_{0.75}}$ (%) | qds $_{\epsilon_{0.5}}$ (%) |
|------------|-------|--------------------------------|---------------------------------|--------------------------------|
| 20ng | 3136 | -4.50 | -14.89 | -31.89 |
| adult | 371 | 3.23 | 7.01 | -13.75 |
| book | 2196 | -10.66 | -19.17 | -40.30 |
| covertypes | 337 | -0.89 | -11.28 | -37.69 |
| kddcup | 285 | -5.26 | -18.95 | -38.95 |
| msnbc | 102 | -7.84 | -33.33 | -63.73 |
| msweb | 1,264 | -2.53 | -3.56 | -34.97 |
| plants | 320 | -6.25 | -18.44 | -42.50 |
| r52 | 2713 | -3.65 | -9.14 | -25.14 |
| uscensus | 220 | -10.45 | -20.45 | -37.73 |
| andes | 336 | -0.89 | -7.14 | -22.92 |
| hailfinder | 64 | -1.56 | +6.25 | -15.62 |
| hepar2 | 92 | -3.26 | -21.74 | -30.43 |
| link | 1,146 | -1.83 | -0.44 | -22.43 |
| munin1 | 208 | 0.00 | +0.96 | -9.62 |
| munin2 | 879 | 0.00 | 0.00 | -13.31 |
| munin3 | 898 | 0.00 | 0.00 | -7.80 |
| munin4 | 903 | 0.00 | 0.00 | -8.53 |
| pathfinder | 161 | -4.35 | -8.70 | -24.22 |
| win95pts | 115 | 0.00 | -0.87 | -12.17 |

tween the different variables. For instance, we see that visits of the website’s pages corresponding to nodes in white (e.g. ‘weather’ (V8), ‘health’ (V9) or ‘business’ (V11)) are importantly influenced by whether the user has also visited the frontpage (V1). For example, learnt parameters show that a user who did not visit the website’s frontpage (V1) is about 10 times more likely to have visited the website’s ‘summary’ page (V13) than a user who did visit the frontpage. Such information is much harder to read from the graph learnt with *sota-BNSL* displayed in Figure 1.

6 DISCUSSION

We have seen that, both in theory and in practice, the quasi-determinism screening approach enables a decrease in computational time and complexity for a small decrease in graph score. This tradeoff is all the more advantageous that there actually are strong pairwise relationships in the data, that can be detected during the screening phase, thus enabling a decrease in the number of variables to be considered by the state of the art structure learning algorithm during the second phase of Algorithm 2.

Optimal cases for this algorithm take place when $n_r(\epsilon)$ is significantly smaller than n for ϵ reasonably small compared to the variable’s entropies. In practice this is reasonably frequent (e.g. 20 newsgroup, msnbc, munin2-4, webkb among others).

Experiments presented in this paper were only conducted on open-source datasets that are frequently used by the community. However, we have also tested our algorithm on industrial descriptive metadatasets from the IoT domain, for which many variables possess (empirically) deterministic parents: we have $n_r(\epsilon = 0)$ very small with respect to n (typically $n_r(\epsilon = 0) \approx 10$ for $n \approx 100$). In this context *qds-BNSL* is up to 20 times faster than *sota-BNSL*, with equivalent or even better learned graphs in terms of BDeu score. This is a specific case (data stored in relational databases in almost perfect third normal form), but it is still interesting since an important amount of data accessible today contains determinism, as previously noted.

Besides, we still have potential to improve the *qds-BNSL* algorithm, by parallelizing the computation of \mathbb{H}^D , and implementing it in C instead of R.

Our main research perspective is to be able to anticipate

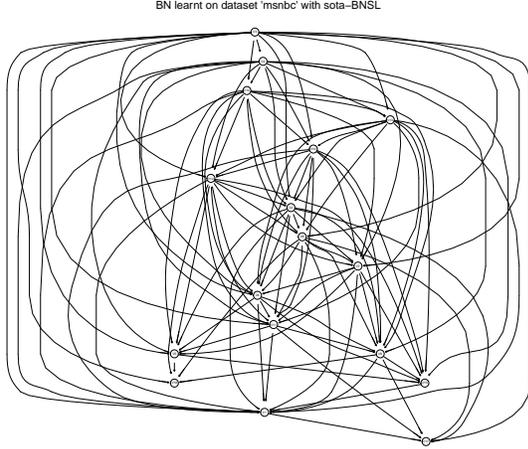


Figure 1: Bayesian network learnt on dataset ‘msnbc’ with *sota-BNSL*. BDeu score: -6.19, Nb of arcs: 102, Time until convergence: 252s

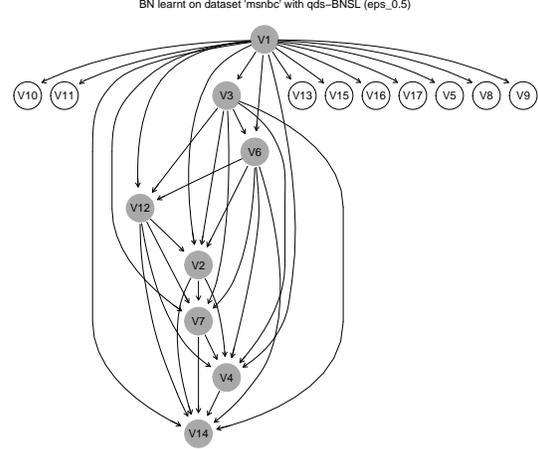


Figure 2: Bayesian network learnt on dataset ‘msnbc’ with *qds-BNSL*($\epsilon_{0.5}$). BDeu score: -6.48, Nb of arcs: 37, Time until convergence: 36s

how good the tradeoff may be before running any algorithm all the way through, saving us from trying *qds-BNSL* on datasets in which there are absolutely no strong pairwise relationships, and enabling us to choose an optimal value of ϵ on datasets for which there is a lot of potential for computational time win with controlled performance loss.

The bound presented in Proposition 4 concerns the MLL score and is far from tight in practice. However, if we could find a tight bound on the BDeu score of the graphs generated by *qds-BNSL*, it would be much easier to estimate the most promising value of ϵ for a given dataset.

Finally, we have some insights on ways to generalize our quasi-determinism screening idea.

The proof of Proposition 2 suggests that the result still holds when F is *any kind of deterministic DAG* (and not only a forest). We could therefore use techniques that detect determinism in a broader sense than only pairwise, to make the screening more efficient. For this purpose we could take inspiration from papers of the knowledge discovery in databases (KDD) community, as Huhtala et al. (1999), or more recently Papenbrock et al. (2015) who evaluate functional dependencies discovery methods.

We also could broaden our definition of quasi-determinism: instead of considering the information-theoretic quantity $H^D(X|Y)$ to describe the strength of the relationship $Y \rightarrow X$, one could choose $\frac{H^D(X|Y)}{H^D(X)}$, which represents the proportion of X ’s entropy that is explained by Y . Moreover, $\frac{H^D(X|Y)}{H^D(X)} \leq \epsilon$ can be rewritten as $\frac{MI^D(X,Y)}{H(X)} \geq 1 - \epsilon$, which gives another insight to quasi-determinism screening: for a given variable X , this comes down to finding a variable Y such

that $MI^D(X, Y)$ is high. This is connected to the idea of Chow and Liu (1968), and later Cheng et al. (1997), for whom pairwise empirical mutual information is central. This alternate definition of ϵ -quasi-determinism does not change the algorithms and complexity considerations described in section 4. Lastly, we could consider other definitions of entropy as the ones presented by Rényi et al. (1961).

References

- Bouckaert, R. (1995). *Bayesian belief networks: from inference to construction*. PhD thesis, Faculteit Wiskunde en Informatica, Utrecht University.
- Buntine, W. (1991). Theory refinement on Bayesian networks. In *Uncertainty Proceedings 1991*, pages 52–60. Elsevier.
- Chen, X.-W., Anantha, G., and Lin, X. (2008). Improving Bayesian network structure learning with mutual information-based node ordering in the K2 algorithm. *IEEE Transactions on Knowledge and Data Engineering*, 20(5):628–640.
- Cheng, J., Bell, D. A., and Liu, W. (1997). Learning belief networks from data: An information theory based approach. In *Proceedings of the sixth international conference on Information and knowledge management*, pages 325–331. ACM.
- Chickering, D. M. (1996). Learning Bayesian networks is NP-complete. *Learning from data: Artificial intelligence and statistics V*, 112:121–130.
- Chow, C. and Liu, C. (1968). Approximating discrete probability distributions with dependence trees. *IEEE transactions on Information Theory*, 14(3):462–467.

- Cooper, G. F. and Herskovits, E. (1992). A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, 9(4):309–347.
- Davis, J. and Domingos, P. (2010). Bottom-up learning of Markov network structure. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 271–278.
- de Campos, C. P. d. and Ji, Q. (2011). Efficient structure learning of Bayesian networks using constraints. *Journal of Machine Learning Research*, 12(Mar):663–689.
- de Morais, S. R., Aussem, A., and Corbex, M. (2008). Handling almost-deterministic relationships in constraint-based Bayesian network discovery: Application to cancer risk factor identification. In *European Symposium on Artificial Neural Networks, ESANN’08*.
- Dheeru, D. and Karra Taniskidou, E. (2017). UCI machine learning repository.
- El Kaed, C., Leida, B., and Gray, T. (2016). Building management insights driven by a multi-system semantic representation approach. In *Internet of Things (WF-IoT), 2016 IEEE 3rd World Forum on*, pages 520–525. IEEE.
- Friedman, N., Nachman, I., and Peér, D. (1999). Learning Bayesian network structure from massive datasets: the ‘sparse candidate’ algorithm. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pages 206–215. Morgan Kaufmann Publishers Inc.
- Heckerman, D., Geiger, D., and Chickering, D. M. (1995). Learning Bayesian networks: The combination of knowledge and statistical data. *Machine Learning*, 20(3):197–243.
- Huhtala, Y., Kärkkäinen, J., Porkka, P., and Toivonen, H. (1999). Tane: An efficient algorithm for discovering functional and approximate dependencies. *The computer journal*, 42(2):100–111.
- Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT press.
- Koo, D. D., Lee, J. J., Sebastiani, A., and Kim, J. (2016). An internet-of-things (iot) system development and implementation for bathroom safety enhancement. *Procedia Engineering*, 145:396–403.
- Luo, W. (2006). Learning Bayesian networks in semi-deterministic systems. In *Canadian Conference on AI*, pages 230–241. Springer.
- Mabrouk, A., Gonzales, C., Jabet-Chevalier, K., and Chojnacki, E. (2014). An efficient Bayesian network structure learning algorithm in the presence of deterministic relations. In *Proceedings of the Twenty-first European Conference on Artificial Intelligence*, pages 567–572. IOS Press.
- Nie, S., de Campos, C. P., and Ji, Q. (2016). Learning Bayesian networks with bounded tree-width via guided search. In *AAAI*, pages 3294–3300.
- Papenbrock, T., Ehrlich, J., Marten, J., Neubert, T., Rudolph, J.-P., Schönberg, M., Zwiener, J., and Naumann, F. (2015). Functional dependency discovery: An experimental evaluation of seven algorithms. *Proceedings of the VLDB Endowment*, 8(10):1082–1093.
- Rényi, A. et al. (1961). On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*. The Regents of the University of California.
- Scanagatta, M., Corani, G., de Campos, C. P., and Zaffalon, M. (2016). Learning treewidth-bounded Bayesian networks with thousands of variables. In *Advances in Neural Information Processing Systems*, pages 1462–1470.
- Scanagatta, M., de Campos, C. P., Corani, G., and Zaffalon, M. (2015). Learning Bayesian networks with thousands of variables. In *Advances in Neural Information Processing Systems*, pages 1864–1872.
- Schwarz, G. et al. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2):461–464.
- Scutari, M. (2009). Learning Bayesian networks with the bnlearn R package. *arXiv preprint arXiv:0908.3817*.
- Silander, T. and Myllymaki, P. (2012). A simple approach for finding the globally optimal Bayesian network structure. *arXiv preprint arXiv:1206.6875*.
- Spirtes, P., Glymour, C. N., and Scheines, R. (2000). *Causation, prediction, and search*. MIT press.
- Teyssier, M. and Koller, D. (2012). Ordering-based search: A simple and effective algorithm for learning Bayesian networks. *Proceedings of the 28th conference on Uncertainty in artificial intelligence*.
- Tsamardinos, I., Brown, L. E., and Aliferis, C. F. (2006). The max-min hill-climbing Bayesian network structure learning algorithm. *Machine Learning*, 65(1):31–78.
- Yaramakala, S. and Margaritis, D. (2005). Speculative Markov blanket discovery for optimal feature selection. In *Fifth IEEE international conference on Data Mining*, pages 4–pp. IEEE.
- Yuan, C., Malone, B., et al. (2013). Learning optimal Bayesian networks: A shortest path perspective. *Journal of Artificial Intelligence Research*.

7 Supplementary material

Proof of Lemma 1: First let us rewrite the MLL score in terms of data counts. We denote $x_i[m]$ the m^{th} observation of variable X_i in the dataset D . For a given $G \in \text{DAG}_V$ and $\theta \in \Theta_G$,

$$\begin{aligned} l(\theta : D) &= \sum_{m=1}^M \log \left(\underbrace{p_\theta(x_1[m], \dots, x_n[m])}_{\prod_{i=1}^n \theta_{x_i[m]|\mathbf{x}_{\pi(i)}[m]}} \right) \\ &= \sum_{m=1}^M \sum_{i=1}^n \log(\theta_{x_i[m]|\mathbf{x}_{\pi(i)}[m]}) \\ &= \sum_{i=1}^n \sum_{x_i, \mathbf{x}_{\pi(i)}} C^D(x_i, \mathbf{x}_{\pi(i)}) \log(\theta_{x_i|\mathbf{x}_{\pi(i)}}) \end{aligned}$$

where $C^D(\cdot)$ is the count function associated with D :

$$\forall I \subset V, C^D(\mathbf{x}_I) = \sum_{m=1}^M \mathbb{I}_{\mathbf{x}_I[m]=\mathbf{x}_I} = Mp^D(\mathbf{x}_I).$$

Moreover, it is well known that for categorical variables, the maximum likelihood estimator θ^{MLE} is given by the local empirical frequencies *i.e.*

$$\theta_{x_i|\mathbf{x}_{\pi(i)}}^{MLE} = p^D(x_i|\mathbf{x}_{\pi(i)}) = \frac{C^D(x_i, \mathbf{x}_{\pi(i)})}{C^D(\mathbf{x}_{\pi(i)})}.$$

Therefore we get:

$$\begin{aligned} s^L(G : D) &= \max_{\theta \in \Theta_G} l(\theta : D) \\ &= l(\theta^{MLE} : D) \\ &= \sum_{i=1}^n \sum_{x_i, \mathbf{x}_{\pi(i)}} C^D(x_i, \mathbf{x}_{\pi(i)}) \log(\theta_{x_i|\mathbf{x}_{\pi(i)}}^{MLE}) \\ &= \sum_{i=1}^n \sum_{x_i, \mathbf{x}_{\pi(i)}} Mp^D(x_i, \mathbf{x}_{\pi(i)}) \log(p^D(x_i|\mathbf{x}_{\pi(i)})) \\ &= -M \sum_{i=1}^n H^D(X_i|\mathbf{X}_{\pi(i)}). \end{aligned}$$

■

Proof of Proposition 1: Let $G \in \text{DAG}_V$ with $V = \llbracket 1, n \rrbracket$ and D containing observations of X_1, \dots, X_n respectively associated with nodes in V .

First, we notice that $s^L(G : D)$ is upper-bounded by the score of a dense DAG. We have shown in Lemma 1 that:

$$s^L(G : D) = -M \sum_{i=1}^n H^D(X_i|\mathbf{X}_{\pi(i)}).$$

It is commonly known that all DAGs are compatible with at least one ordering of the nodes, *i.e.* that $\exists \sigma \in \mathcal{S}_n$ such that

$$\forall i, j \in V \text{ s.t. } j \in \pi^G(i), \sigma(j) < \sigma(i).$$

In other words, σ represents an ordering in which each node comes after its parents.

Let $\sigma \in \mathcal{S}_n$ be an ordering compatible with G . Using the fact that for any variables X, Y, Z , $H^D(X|Y) \geq H^D(X|Y, Z)$ we then get that $\forall i \in V \setminus \{\sigma^{-1}(1)\}$,

$$H^D(X_i|\mathbf{X}_{\pi(i)}) \geq H^D(X_i|\mathbf{X}_{\sigma^{-1}(\{1, \dots, \sigma(i)-1\})}).$$

Plugging this inequality in the first equation, reordering the sum according to σ , and using the chain rule for entropies, we get:

$$\begin{aligned} -\frac{s^L(G : D)}{M} &\geq \sum_{i=1}^n H^D(X_i|\mathbf{X}_{\sigma^{-1}(\{1, \dots, \sigma(i)-1\})}) \\ &= H^D(X_{\sigma^{-1}(1)}) \\ &\quad + \sum_{\sigma(i)=2}^n H^D(X_{\sigma^{-1}(\sigma(i))}|\mathbf{X}_{\sigma^{-1}(\{1, \dots, \sigma(i)-1\})}) \\ &= H^D(X_{\sigma^{-1}(1)}) \\ &\quad + H^D(X_{\sigma^{-1}(2)}|X_{\sigma^{-1}(1)}) + \dots \\ &\quad + H^D(X_{\sigma^{-1}(n)}|X_{\sigma^{-1}(1)}, \dots, X_{\sigma^{-1}(n-1)}) \\ &= H^D(X_{\sigma^{-1}(1)}, \dots, X_{\sigma^{-1}(n)}) \\ &= H^D(X_1, \dots, X_n), \end{aligned}$$

which gives

$$s^L(G : D) \leq -M H^D(X_1, \dots, X_n).$$

Let T be as in the hypothesis of Proposition 1, we are now going to prove that this bound is reached for T which will give us the wanted result.

Without any loss of generality, let us suppose that T 's root is 1. Then,

$$\begin{aligned} s^L(T : D) &= -M \sum_{i=1}^n H^D(X_i|\mathbf{X}_{\pi(i)}) \\ &= -M \left(H^D(X_1) + \underbrace{\sum_{i=2}^n H^D(X_i|\mathbf{X}_{\pi(i)})}_{=0} \right) \\ &\geq -M H^D(X_1, \dots, X_n) \\ &= \max_{G \in \text{DAG}_V} s^L(G : D). \end{aligned}$$

■

Proof of Proposition 2: Let $F = \bigcup_{k=1}^p T_k$ and $G_{R(F)}^*$

be as in the Proposition's hypotheses. Without loss of generality, we consider i to be the root of the tree T_i . Therefore, $R(F) = \llbracket 1, p \rrbracket$.

Let us also define the following *root* function that associates to each node the root of the tree it belongs to:

$$r : \begin{cases} V & \longrightarrow R(F) \\ i & \longmapsto k \text{ s.t. } X_i \in V_{T_k}. \end{cases}$$

Let $G_{R(F)}^* \in DAG_{R(F)}$ such that:

$$G_{R(F)}^* \in \operatorname{argmax}_{G \in DAG_{R(F)}} s^L(G : D)$$

and $G^* = F \cup G_{R(F)}^*$ i.e.

- $V_{G^*} = V$
- $A_{G^*} = (\bigcup_{k=1}^p A_{T_k}) \cup A_{G_{R(F)}^*}$

We will show as in the proof of Proposition 1 that

$$s_{DAG_V}^L(G^* : D) \geq \max_{G \in DAG_V} s_{DAG_V}^L(G : D)$$

which implies that $G^* \in \operatorname{argmax}_{G \in DAG_V} s_{DAG_V}^L(G : D)$.

We write:

$$\begin{aligned} s_{DAG_V}^L(G^*) &= -M \sum_{i=1}^n H^D(X_i | \mathbf{X}_{\pi_{G^*}(i)}) \\ &= \underbrace{-M \sum_{i=1}^p H^D(X_i | \mathbf{X}_{\pi_{G^*}(i)})}_{(a)} \\ &\quad \underbrace{-M \sum_{i=p+1}^n H^D(X_i | \mathbf{X}_{\pi_{G^*}(i)})}_{(b)} \end{aligned}$$

We then compute separately the terms (a) and (b):

- *Computation of (a)*

The first term corresponds to the score of the graph

$G_{R(F)}^*$ as an element of $DAG_{R(F)}$.

Indeed, by construction of G^* ,

$$\forall i \in R(F), \pi^{G^*}(i) = \pi^{G_{R(F)}^*}(i).$$

Moreover, $G_{R(F)}^*$ maximizes the MLL score on $DAG_{R(F)}$. We can now write:

$$\begin{aligned} (a) &= -M \sum_{i=1}^p H^D(X_i | \mathbf{X}_{\pi_{G^*}(i)}) \\ &= -M \sum_{i=1}^p H^D(X_i | \mathbf{X}_{\pi_{G_{R(F)}^*}(i)}) \\ &= s^L(G_{R(F)}^* : D) \\ &= \max_{G \in DAG_{R(F)}} s^L(G : D_{R(F)}) \\ &= -MH^D(X_1, \dots, X_p). \end{aligned}$$

- *Computation of (b)*

By construction of G^* ,

$$\forall i \in V \setminus R(F), \pi^{G^*}(i) = \pi^{T_{r(i)}}(i).$$

Moreover since the T_k 's are deterministic trees, it follows that

$$\forall i \in V \setminus R(F), H^D(X_i | \mathbf{X}_{\pi^{T_{r(i)}}(i)}) = 0.$$

Therefore we can write

$$\begin{aligned} (b) &= -M \sum_{i=p+1}^n H^D(X_i | \mathbf{X}_{\pi_{G^*}(i)}) \\ &= -M \sum_{i=p+1}^n H^D(X_i | \mathbf{X}_{\pi^{T_{r(i)}}(i)}) \\ &= 0. \end{aligned}$$

Collecting the above results yields

$$\begin{aligned} s_{DAG_V}^L(G^*) &= (a) \\ &= -MH^D(X_1, \dots, X_p) \\ &\geq -MH^D(X_1, \dots, X_n) \\ &= \max_{G \in DAG_V} s^L(G : D). \end{aligned}$$

■

In a few words: the idea of the proof relies on the fact that $H^D(\mathbf{X}) = H^D(\mathbf{X}_{R(F)})$: the information associated with all the variables X_1, \dots, X_n is entirely contained in the root variables $\mathbf{X}_{R(F)}$.

Proof of Proposition 3: Let D and ϵ be objects that satisfy the input constraints of Algorithm 1, and let F_ϵ be the object that is returned by Algorithm 1 with inputs D and ϵ .

F_ϵ is a **directed graph**, by definition. Moreover, it is built so that all of its nodes had at most one parent (line 12).

To conclude, we therefore only have to prove that F_ϵ does not contain cycles.

Let us suppose that there is a cycle i_1, \dots, i_p in F_ϵ . There are two cases:

1. Either all associated variables have the same entropy: $H^D(X_{i_1}) = H^D(X_{i_2}) = \dots = H^D(X_{i_p})$. In which case, there is necessarily two successive nodes in the cycle i_l, i_{l+1} such that $i_l < i_{l+1}$. However, $H^D(X_{i_l} | X_{i_{l+1}}) = H^D(X_{i_{l+1}} | X_{i_l}) \leq \epsilon$, which means that when the algorithm reaches line 4, $i_{l+1} \in \pi_\epsilon(i_l)$ and $i_l \in \pi_\epsilon(i_{l+1})$. Since i_l is treated before i_{l+1} in the for loop, this would result in i_l being removed from $\pi_\epsilon(i_{l+1})$, thus preventing for i_l to ever be i_{l+1} 's parent: we have a **contradiction**.
2. Either there exist at least two variables in the cycle that do not have the same entropy: $H^D(X_{i_k}) \neq H^D(X_{i_{k'}})$, for $k, k' \in \llbracket 1, p \rrbracket$.

In this case, there also exist two successive nodes i_l, i_{l+1} such that $H^D(X_{i_l}) \neq H^D(X_{i_{l+1}})$. Let us suppose that $H^D(X_{i_{l+1}}) > H^D(X_{i_l})$. In that case:

$$\begin{aligned} H^D(X_{i_l}|X_{i_{l+1}}) &= H^D(X_{i_{l+1}}|X_{i_l}) \\ &\quad + \underbrace{H^D(X_{i_l}) - H^D(X_{i_{l+1}})}_{<0} \\ &< H^D(X_{i_{l+1}}|X_{i_l}) \\ &\leq \epsilon. \end{aligned}$$

Therefore when the algorithm reaches line 4, $i_{l+1} \in \pi_\epsilon(i_l)$. But when treating either i_l or i_{l+1} during the for loop of lines 4-9, the test on line 6 necessarily implies that i_l is removed from $\pi_\epsilon(i_{l+1})$ (since $H^D(X_{i_l}|X_{i_{l+1}}) < H^D(X_{i_{l+1}}|X_{i_l})$). This is in **contradiction** with the fact that i_l is i_{l+1} 's parent.

Therefore, $H^D(X_{i_{l+1}}) < H^D(X_{i_l})$, and arcs in A_{F_ϵ} follow *nonincreasing entropies*, with at least one decrease since we suppose that all entropies are not equal. This is not possible in a cycle: there is a **contradiction**.

We conclude that there is **no cycle** in F_ϵ . ■

Proof of Proposition 4: The structure of the proof is the same as the one from Proposition 2. The only difference lies in the computation of term (b):

$$\begin{aligned} (b) &= -M \sum_{i=p+1}^n H^D(X_i|\mathbf{X}_{\pi_{G^*}(i)}) \\ &= -M \sum_{i=p+1}^n \underbrace{H^D(X_i|\mathbf{X}_{\pi_{T_r(i)}(i)})}_{\leq \epsilon} \\ &\geq -M(n-p)\epsilon \\ &\geq -Mn\epsilon. \end{aligned}$$

plugging this in the separated expression of the MLL score of G^* in terms (a) and (b) yields the wanted result. ■