



Inversion of Bayesian networks

Jesse van Oostrum^{a,*}, Peter van Hintum^b, Nihat Ay^{a,c,d}

^a Institute for Data Science Foundations, Hamburg University of Technology, Hamburg, Germany

^b New College, University of Oxford, Oxford, UK

^c Leipzig University, Leipzig, Germany

^d Santa Fe Institute, Santa Fe, USA

ARTICLE INFO

Keywords:

Graphical models
Variational inference
Amortized inference
Bayesian networks
Variational autoencoder
Generative model
Recognition model

ABSTRACT

Variational autoencoders and Helmholtz machines use a recognition network (encoder) to approximate the posterior distribution of a generative model (decoder). In this paper we establish some necessary and some sufficient properties of a recognition network so that it can model the true posterior distribution exactly. These results are derived in the general context of probabilistic graphical modelling / Bayesian networks, for which the network represents a set of conditional independence statements. We derive both global conditions, in terms of d-separation, and local conditions for the recognition network to have the desired qualities. It turns out that for the local conditions the *perfectness* property (for every node, all parents are joined) plays an important role.

1. Introduction

A generative model is a set of probability distributions that models the distribution of observed and latent variables. Generative models are used in many machine learning applications. One is often interested in inferring the latent variable on the basis of an observation, i.e. obtaining the posterior distribution. For complex generative models it is often hard to calculate the posterior distribution analytically. The field of variational Bayesian inference [18] studies different ways of approximating the true posterior. One approach within this field is called *amortised inference* [7]. This approach distinguishes itself through using one set of parameters for recognition that is optimised over multiple data points. This can be contrasted with “memoryless” inference algorithms, such as the message passing algorithm [15,3], which finds a separate set of parameters for every data point. Both the variational autoencoder (VAE) [8] and the Helmholtz machine [4] are examples of amortised inference. In their most general form, these consist of a Bayesian network that is used to model the generative distribution. A second network, called the recognition model, is used to model the posterior distribution. Both these networks have the same set of nodes, namely the union of the observed and latent variables. In the generative network the arrows point from the latent to the observed nodes but in the recognition network it is the other way around. The recognition network is therefore in some sense an inversion of the generative network. In many applications, one simply reverses the direction of the edges of the generative network to obtain the recognition network. However, as the simple example in Fig. 1 shows, this does not guarantee that the recognition model is actually able to model the true posterior distribution of the generative model. In this paper, we establish some necessary and some sufficient properties of the recognition network such that we do have this guarantee. We first discuss these properties in terms of d-separation, subsequently in terms of perfectness, and finally in terms of single edge operations using the Meek conjecture [13].

* Corresponding author.

E-mail address: jesse.van@tuhh.de (J. van Oostrum).

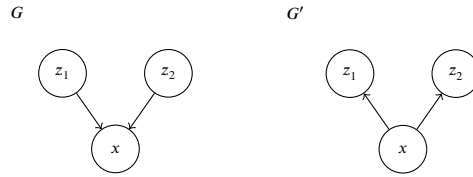


Fig. 1. Pair of DAGs G (left) and G' (right) where G' is obtained by reversing the direction of the edges in G . The variables z_1, z_2 represent the latent variables and x the observed variable. The distribution p such that z_1, z_2 are independent Bernoulli(0.5) and $x = z_1 + z_2 \bmod 2$ can be modelled by G , but the conditional distribution $p_{z_1, z_2 | x}$ cannot be modelled by G' .

In practice, one often puts further restrictions on the probability distributions the networks can model. For example, the distribution of an individual node can be required to be Gaussian. The mean (and variance) are in this case a function of the values of the parent nodes. We discuss the case of a restricted set of probability distributions in the last part of the results section.

The question of finding a sparse G' that can approximate the posterior distribution of the generative model well is also studied from a more practical perspective, using methods from machine learning. One can use a sparsity prior when determining the recognition model, to encourage that only the edges really necessary for modelling the posterior are added. Löwe et al. [12], Louizos et al. [11], Molchanov et al. [14] present several approaches.

Markov equivalence is a property of a pair of Bayesian networks that indicates that they encode the same set of conditional independence statements [17,6]. A generalisation of this, which we will call *Markov inclusion*, is when the set of conditional independence statements encoded in one graph is a subset of the conditional independence statements encoded in the other graph [1]. We will see in Proposition 1 that the results in this paper can also be viewed as describing under which conditions one Bayesian network is Markov inclusive of another.

Webb et al. [19] deal with a closely related problem. They present an algorithm for inverting the generative network that gives a recognition network with the desired properties. While the present article also deals with the algorithmic aspects, it puts more emphasis on the conditions for the recognition network to have the desired properties. The authors were unaware of the publication [19] prior to the acceptance of the present paper.

1.1. Example

Before giving a formal definition of the problem, we illustrate the topic of this paper by an intuitive example that provides context for the rest of the paper. Consider the generative model for diseases and their symptoms in Fig. 2. Our goal is to find a model to

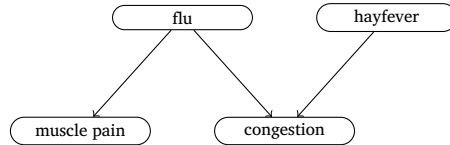


Fig. 2. Example from [9].

perform inference on this generative model, i.e. to find the posterior distribution $P(\text{diseases} \mid \text{symptoms})$. Naively, one could reverse the edges of the generative model to obtain the recognition model in Fig. 3.

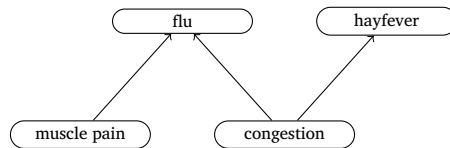


Fig. 3. Recognition model obtained from reversing the edges of the generative model.

It is clear that when someone is congested, whether or not they have muscle pain affects the likelihood of that person having hayfever. If someone is congested and also has muscle pain, the congestion is more likely to be caused by the flu. On the other hand, the absence of muscle pain would make it more likely for the congestion to be caused by hayfever. This dependence is however not captured in the graph in Fig. 3, because no information can flow from muscle pain to hayfever. By adding an edge between muscle pain and hayfever, or between flu and hayfever, this dependence can be captured. (Fig. 4.)

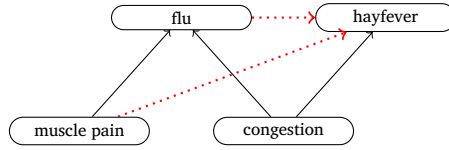


Fig. 4. Recognition model with optional arrows in red to capture the dependence between muscle pain and hayfever.

2. Notation

2.1. Graph theory

For a comprehensive overview of the theory and terminology of probabilistic graphical models, we refer to [10,3,16]. We define a *graph* G to be a pair $G = (N, E)$ where N is a non-empty finite set of *vertices* or *nodes* and $E \subset N \times N$ such that $(s, s) \notin E$ for all $s \in N$. A graph $H = (A, \tilde{E})$ is called a *subgraph* of G if $A \subset N$ and $\tilde{E} \subset E$ and we write $H \subset G$. For a subset $A \subset N$, the *vertex-induced subgraph* of G is denoted $G[A]$ and is given by $(A, E[A])$, with $E[A] = \{(s, t) \in E : s, t \in A\}$. When both (s, t) and (t, s) are in E , we say that there is an *undirected edge* between s and t . When $(s, t) \in E$ and $(t, s) \notin E$ we say that there is a *directed edge* going from s to t and write $s \rightarrow t$. We say that a graph is *directed* if all edges are directed.

In the following we let $G = (N, E)$ be a directed graph. We say that two vertices $s, t \in N$ are *joined* if there is an edge between the two. A set of vertices is called *complete* if all pairs of its elements are joined. A *path* in G from s to t is a sequence of distinct nodes $s = u_0, \dots, u_n = t$ such that $(u_i, u_{i+1}) \in E$ for all $i \in \{0, \dots, n-1\}$. A *cycle* is a path of length $n > 1$ with the modification that the end points are identical. We say that a graph is *acyclic* if it does not possess any cycles. A directed graph which is acyclic is called a *directed acyclic graph*, or DAG.

In the following we let $G = (N, E)$ be a DAG. A *trail* from s to t is a sequence of distinct nodes $s = u_0, \dots, u_n = t$ such that $(u_i, u_{i+1}) \in E$ or $(u_{i+1}, u_i) \in E$ or both, for all $i \in \{0, \dots, n-1\}$. Note that movement along a trail could go against the direction of the arrows, in contrast to the case of a path. A *loop* is a trail of length $n > 1$ with the modification that the end points are identical. If $(s, t) \in E$ we call s a *parent* of t and t a *child* of s . The set of parents of a node t is denoted $\text{pa}_G(t)$ and the set of children of a node s is denoted $\text{ch}_G(s)$. A node $t \neq s$ is called a *descendant* of s if there exists a path from s to t . The set of descendants of a node s is denoted $\text{des}_G(s)$. The set of *non-descendants* of s is given by $N \setminus (\{s\} \cup \text{des}_G(s))$ and is denoted by $\text{nd}_G(s)$. G is called *perfect* if for all s , the set $\text{pa}_G(s)$ is complete. We let $\text{Leaves}(G) = \{s \in N : \text{ch}_G(s) = \emptyset\}$ be the set of nodes without children, and $\text{Roots}(G) = \{s \in N : \text{pa}_G(s) = \emptyset\}$ be the set of nodes without parents (see Fig. 5). For $e = (s, t) \in E$, let $e^* = (t, s)$, $E^* = \{e^* : e \in E\}$, $G^* = (N, E^*)$ the graph G with its edges reversed, $E^\sim = E \cup E^*$, $G^\sim = (N, E^\sim)$ the *skeleton* (i.e. undirected version) of G , and $E_{\text{pa}(G)}^\sim = \{(t_1, t_2) : \exists s \in N, t_1, t_2 \in \text{pa}_G(s)\}$, $G^M = (N, E^\sim \cup E_{\text{pa}(G)}^\sim)$ the *moral graph* of G , which is the skeleton of G , with extra (undirected) edges between all parents of every vertex in G .

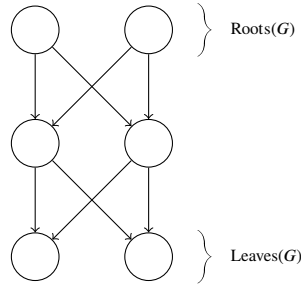
For a trail $\gamma = (u_0, \dots, u_n)$, we let $\gamma^* = (u_n, \dots, u_0)$ be its reversed version, and $(\gamma, s) = (u_0, \dots, u_n, s)$ for $s \notin \gamma$ such that $(u_n, s) \in E$ or $(s, u_n) \in E$ be its prolongation. Now let $\gamma_1 = (u_0, \dots, u_n)$ and $\gamma_2 = (v_0, \dots, v_m)$ be trails such that $u_n = v_0$. We define $\gamma = \gamma_1 \gamma_2$ to be the *concatenation* of γ_1 and γ_2 such that $\gamma = (u_0, \dots, u_i, s, v_j, \dots, v_m)$, with s the first node in γ_1 which belongs to $\gamma_1 \cap \gamma_2$. Let $\gamma = (u_0, \dots, u_n)$ be a trail and u_i a node on this trail that is not one of the endpoints, i.e. $0 < i < n$. u_i is called a *v-structure* if $u_{i-1} \rightarrow u_i \leftarrow u_{i+1}$. γ is said to be *blocked* by $S \subset N$ if γ contains a vertex u such that either: (a) $u \in S$ and u is not a v-structure; (b) u and $\text{des}_G(u)$ are not in S and u is a v-structure. Note that when one of the endpoints of the trail is in S , the trail is definitely blocked, since endpoints can never be v-structures. Let $A, B, S \subset N$ (not necessarily disjoint). A, B are said to be *d-separated* by S if all trails from A to B are blocked by S and we write $A \perp_G B \mid S$. A *topological ordering* of G is an injective map $\mathcal{O} : N \rightarrow \{1, \dots, |N|\}$ that assigns to every node a number such that, if two nodes are joined, the edge points from the lower to the higher numbered node, i.e. $(s, t) \in E$ implies $\mathcal{O}(s) < \mathcal{O}(t)$. When $s, t \in N$ are such that $\mathcal{O}(s) < \mathcal{O}(t)$ we say s is *older* than t , and t is *younger* than s . Given a topological ordering \mathcal{O} , the set of *predecessors* of a node s , denoted $\text{pr}^\mathcal{O}(s)$, is the set of all nodes with a lower topological number,¹ i.e. $\text{pr}^\mathcal{O}(s) = \{t \in N : \mathcal{O}(t) < \mathcal{O}(s)\}$. Note that, although $\text{des}_G(s) \cap \text{pr}^\mathcal{O}(s) = \emptyset$, the set $\text{pr}^\mathcal{O}(s)$ in general depends on the choice of topological ordering (see Fig. 6).

2.2. Probability on graphs

Consider a DAG $G = (N, E)$. To every node $s \in N$ we associate a measurable space (X_s, \mathcal{X}_s) . The state spaces are either real finite-dimensional vector spaces or finite sets and to each measurable space we associate a (σ -finite) base measure μ_s which is typically the Lebesgue measure or counting measure respectively. For a non-empty² subset $A \subset N$ we let $X_A = \times_{s \in A} X_s$ be the Cartesian product of the individual X_s and $\mathcal{X}_A = \otimes_{s \in A} \mathcal{X}_s$ be the product σ -algebra. We write $(X, \mathcal{X}) = (X_N, \mathcal{X}_N)$ and assign to this space the base measure $\mu = \otimes_{s \in N} \mu_s$, which is the product measure. In this paper, we consider probability distributions P over the space (X, \mathcal{X}) . For every $s \in N$ we let $X_s : X \rightarrow X_s$ be the random variable projecting onto the individual spaces, and similarly $X_A = (X_s)_{s \in A}$ and $X = X_N$. A

¹ Note that some authors define the set of predecessors to be the set of nodes with a lower topological number in all possible topological orderings.

² When $A = \emptyset$ we let $X_\emptyset = \{\emptyset\}$ be the space containing only the empty sequence and $\mathcal{X} = \{\emptyset, X_\emptyset\}$ be the trivial σ -algebra. X_\emptyset is simply given by the constant map $X_\emptyset : \mathcal{X} \ni x \mapsto \emptyset \in X_\emptyset$. See [16] Section 2.1 for details.

Fig. 5. Different subsets of N for a graph G .

typical element of X_s is denoted x_s with $x_A = (x_s)_{s \in A}$ and $x = (x_s)_{s \in N}$. We write P_A for the distribution of X_A on (X_A, \mathcal{X}_A) , i.e. for $A \in \mathcal{X}_A$,

$$P_A(A) := P\left((X_A)^{-1}(A)\right) = P(A \times X_{N \setminus A}). \quad (1)$$

For $A, C \subset N$, we say that a map $K : \mathcal{X}_A \times X_C \rightarrow [0, 1]$ is a *Markov kernel* if

$$A \mapsto K(A, x_C) \quad \text{is a probability measure on } \mathcal{X}_A \text{ for all } x_C \in X_C, \quad (2)$$

$$x_C \mapsto K(A, x_C) \quad \text{is } \mathcal{X}_C\text{-measurable for all } A \in \mathcal{X}_A. \quad (3)$$

Furthermore, we say that K is a (*regular*) *version of the conditional probability of A given C* if it is Markov kernel and for all $A \in \mathcal{X}_A, C \in \mathcal{X}_C$

$$P_{A \cup C}(A \times C) = \int_C K(A, x_C) dP_C(x_C) \quad (4)$$

holds. It can be shown that in our setting, one can always find such a Markov kernel that is unique P_C -a.e. [5]. We therefore also denote such a Markov kernel by $P_{A|C}$. For $A, B, C \subseteq N$ we say that A is *conditionally independent* of B given C and write $A \perp\!\!\!\perp B | C$ if for every $A \in \mathcal{X}_A$ and $B \in \mathcal{X}_B$

$$P_{A \cup B|C}(A \times B | x) = P_{A|C}(A | x) \cdot P_{B|C}(B | x) \quad P_C\text{-a.e.} \quad (5)$$

For $s \in N$, a *kernel function* will be a map $k^s(\cdot | \cdot) : X_s \times X_{\text{pa}_G(s)} \rightarrow \mathbb{R}_{\geq 0}$ such that for all $x_{\text{pa}_G(s)} \in X_{\text{pa}_G(s)}$

$$\int k^s(x_s | x_{\text{pa}_G(s)}) d\mu_s(x_s) = 1. \quad (6)$$

A probability distribution P is said to *factorise* over a DAG G if it has a density p w.r.t. the product measure μ and there exist kernel functions $(k^s)_{s \in N}$ such that

$$p(x) = \prod_{s \in N} k^s(x_s | x_{\text{pa}_G(s)}). \quad (7)$$

We denote the set of probability distributions on X that factorise over G by \mathcal{P}^G . Now let $A \subset N$ be such that its parents are themselves in A . A Markov kernel $K : \mathcal{X}_{N \setminus A} \times X_A \rightarrow [0, 1]$ is said to *factorise* over a DAG G if there exist kernel functions $(k^s)_{s \in N \setminus A}$ such that for every $x_A \in X_A$, $K(\cdot, x_A)$ has a density $p(\cdot | x_A)$, such that

$$p(x_{N \setminus A} | x_A) = \prod_{s \in N \setminus A} k^s(x_s | x_{\text{pa}_G(s)}). \quad (8)$$

We denote the set of such Markov kernels by \mathcal{K}^G .

3. Problem statement

Goal I Given a DAG $G = (N, E)$, find a DAG $G' = (N, E')$ such that $\text{Roots}(G') \supset \text{Leaves}(G)$ and for every $P \in \mathcal{P}^G$, there exists a $K \in \mathcal{K}^{G'}$ that is a version of the conditional distribution $P_{N \setminus \text{Leaves}(G) | \text{Leaves}(G)}$.

It turns out (Proposition 1) that this goal is equivalent (up to edges in G' between nodes in $\text{Leaves}(G)$) to the following goal:

Goal II Given a DAG $G = (N, E)$, find a DAG $G' = (N, E')$ such that $\mathcal{P}^{G'} \supset \mathcal{P}^G$ and the parents in G' of the nodes in $\text{Leaves}(G)$ are themselves in $\text{Leaves}(G)$.

Note that even though the parents in G' of the nodes in $\text{Leaves}(G)$ are themselves in $\text{Leaves}(G)$, it is not necessarily the case that $\text{Leaves}(G)$ are oldest in every topological ordering of G' . See Fig. 6 for an example. However, there exists a topological ordering of G' for which the nodes in $\text{Leaves}(G)$ precede the other nodes if and only if the parents in G' of the nodes in $\text{Leaves}(G)$ are themselves in $\text{Leaves}(G)$.

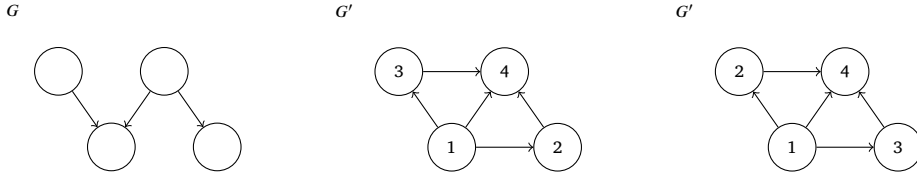


Fig. 6. Pair of DAGs G, G' that satisfy the second requirement of Goal II, but for which there exists a topological ordering of G' (the one on the right) that does not reflect this.

In the remainder of the paper, we will focus on Goal II. Moreover, we sometimes impose the following extra condition:

$$G' \supset G^*. \quad (9)$$

It can be argued that this is a natural condition since it enforces that the hierarchical structure of the generative model G is preserved when finding a suitable G' . Furthermore, note that G^* satisfies the requirement that the parents in G^* of $\text{Leaves}(G)$ are themselves in $\text{Leaves}(G)$.

4. Preliminaries

Lemma 1. Let $G = (N, E_G), H = (N, E_H)$ be DAGs such that $E_G \subset E_H$. Then $\mathcal{P}^G \subset \mathcal{P}^H$.

Proof. Since $\text{pa}_G(s) \subset \text{pa}_H(s)$ for every node s , a density that can be written as $\prod_s k^s(x_s | x_{\text{pa}_G(s)})$ can also be written as $\prod_s k^s(x_s | x_{\text{pa}_H(s)})$. \square

Lemma 2. (Theorem 5.14 in Cowell et al. [3]) Let G be a DAG with a topological ordering \mathcal{O} . For a probability distribution P on \mathbf{X} , the following conditions are equivalent:

- (i) $P \in \mathcal{P}^G$,
- (ii) for all sets $A, B, S \subset N$ such that $A \perp_G B \mid S$ we have $A \perp B \mid S$ w.r.t. P ,
- (iii) for all s we have $s \perp \text{nd}_G(s) \mid \text{pa}_G(s)$ w.r.t. P ,
- (iv) for all s we have $s \perp \text{pr}^{\mathcal{O}}(s) \mid \text{pa}_G(s)$ w.r.t. P .

Corollary 1. Let $\mathcal{O}, \tilde{\mathcal{O}}$ be two topological orderings of G . If P satisfies property (iv) of Lemma 2 w.r.t. \mathcal{O} , then the same is true for $\tilde{\mathcal{O}}$.

Proof. Note that (i) – (iii) of Lemma 2 are independent of the topological ordering. Therefore we have the following implications: for all s we have $s \perp \text{pr}^{\mathcal{O}}(s) \mid \text{pa}_G(s)$ w.r.t. $P \implies P \in \mathcal{P}^G$ (with topological ordering \mathcal{O}) $\implies P \in \mathcal{P}^G$ (with topological ordering $\tilde{\mathcal{O}}$) \implies for all s we have $s \perp \text{pr}^{\tilde{\mathcal{O}}}(s) \mid \text{pa}_G(s)$ w.r.t. P . \square

5. Results

5.1. Equivalence of two goals

Proposition 1. Let G be a DAG. The following statements hold:

- (a) Let $G' = (N, E')$ be a DAG that satisfies Goal I and let \mathcal{O}' be a topological ordering of G' . Then $\tilde{G}' = (N, E' \cup E_{\text{Leaves}})$, with $E_{\text{Leaves}} = \{(s, t) : s, t \in \text{Leaves}(G), \mathcal{O}'(s) < \mathcal{O}'(t)\}$ satisfies Goal II.
- (b) Let $G' = (N, E')$ be a DAG that satisfies Goal II. Then $\tilde{G}' = (N, E' \setminus E_{\text{Leaves}})$, with $E_{\text{Leaves}} = \{(s, t) \in E' : s, t \in \text{Leaves}(G)\}$ satisfies Goal I.

The proof of this proposition is based on the following lemma, in which G plays the role of G' in the proposition.

Lemma 3. Let P be a distribution on \mathbf{X} that has a density p w.r.t. μ .

- (a) Let $G = (N, E)$ be a DAG with topological ordering \mathcal{O} , $A \subset \text{Roots}(G)$, and $H = (N, E \cup E_A)$, with $E_A = \{(s, t) : s, t \in A, \mathcal{O}(s) < \mathcal{O}(t)\}$. If there exists a Markov kernel $K \in \mathcal{K}^G$ that is a version of the conditional distribution $P_{N \setminus A | A}$, then $P \in \mathcal{P}^H$.
- (b) Let $G = (N, E)$ be a DAG, $A \subset N$ such that the parents of nodes in A are themselves in A , and $H = (N, E \setminus E_A)$, $E_A = \{(s, t) \in E : s, t \in A\}$. If $P \in \mathcal{P}^G$, then there exists a Markov kernel $K \in \mathcal{K}^H$ that is a version of the conditional distribution $P_{N \setminus A | A}$.

Proof. (a) Suppose K is a version of $P_{N \setminus A|A}$ and $K \in \mathcal{K}^G$. We need to show $P \in \mathcal{P}^H$. We can write p as follows:

$$p(x) = p(x_{N \setminus A} | x_A) p(x_A), \quad (10)$$

where $p(x_{N \setminus A} | x_A)$ is the density corresponding to K [5]. From the fact that $K \in \mathcal{K}^G$ we know

$$p(x_{N \setminus A} | x_A) = \prod_{s \in N \setminus A} k^s(x_s | x_{\text{pa}_G(s)}). \quad (11)$$

Since all the nodes in A are joined in H we have

$$p(x_A) = \prod_{s \in A} k^s(x_s | x_{\text{pa}_H(s)}). \quad (12)$$

Combining the above gives

$$p(x) = \prod_{s \in N \setminus A} k^s(x_s | x_{\text{pa}_G(s)}) \prod_{s \in A} k^s(x_s | x_{\text{pa}_H(s)}) \quad (13)$$

$$= \prod_{s \in N} k^s(x_s | x_{\text{pa}_H(s)}), \quad (14)$$

and therefore $P \in \mathcal{P}^H$.

(b) Now suppose $P \in \mathcal{P}^G$ and $x \in \mathbf{X}$ such that $p(x_A) > 0$. We can write

$$p(x) = \prod_{s \in N} k^s(x_s | x_{\text{pa}_G(s)}) \quad (15)$$

$$= \prod_{s \in N \setminus A} k^s(x_s | x_{\text{pa}_G(s)}) \prod_{s \in A} k^s(x_s | x_{\text{pa}_G(s)}) \quad (16)$$

$$= \prod_{s \in N \setminus A} k^s(x_s | x_{\text{pa}_H(s)}) \prod_{s \in A} k^s(x_s | x_{\text{pa}_G(s)}), \quad (17)$$

where we can switch from pa_G to pa_H in the third equality because there are only edges removed between nodes in A to obtain H . From the definition of A it follows that $\text{pa}_G(s) \subset A$ for all $s \in A$, hence, $\prod_{s \in A} k^s(x_s | x_{\text{pa}_G(s)}) = p(x_A)$. Dividing by $p(x_A)$ on both sides gives

$$p(x_{N \setminus A} | x_A) = \prod_{s \in N \setminus A} k^s(x_s | x_{\text{pa}_H(s)}). \quad (18)$$

We know that there exists a Markov kernel K that is a version of the conditional distribution of $N \setminus A$ given A and that this kernel has density $p(x_{N \setminus A} | x_A)$ [5]. Equation (18) shows that the density factorises and therefore $K \in \mathcal{K}^H$. \square

5.2. Conditions in terms of d -separation

Necessary and sufficient conditions for our goal can be deduced from the following theorem:

Theorem 1. Let $G = (N, E), G' = (N, E')$ be DAGs, with \mathcal{O}' a topological ordering for G' . The following statements are equivalent:

- (i) $\mathcal{P}^{G'} \supset \mathcal{P}^G$,
- (ii) For all sets $A, B, S \subset N$ such that $A \perp_{G'} B | S$, we have $A \perp_G B | S$,
- (iii) For all $s \in N$, we have $s \perp_G \text{nd}_{G'}(s) | \text{pa}_{G'}(s)$,
- (iv) For all $s \in N$, we have $s \perp_G \text{pr}^{\mathcal{O}'}(s) | \text{pa}_{G'}(s)$.

Proof. (i) \implies (ii) (by contradiction) Suppose there exist A, B, S such that $A \perp_{G'} B | S$, but $A \not\perp_G B | S$. From equation (5.7) in Cowell et al. [3] we know that there exists a $P \in \mathcal{P}^G$ for which $A \not\perp B | S$. This violates (ii) of Lemma 2 and therefore $P \notin \mathcal{P}^{G'}$.

(ii) \implies (i) Let $P \in \mathcal{P}^G$. We need to show $P \in \mathcal{P}^{G'}$. From (i) \implies (ii) in Lemma 2 we know that $A \perp_G B | S \implies A \perp\!\!\!\perp B | S$ for P . Combining this with the assumption $A \perp_{G'} B | S \implies A \perp_G B | S$ gives $A \perp_{G'} B | S \implies A \perp\!\!\!\perp B | S$. This means that P satisfies (ii) of Lemma 2 w.r.t. G' and therefore $P \in \mathcal{P}^{G'}$.

(i) \iff (iii) and (i) \iff (iv) can be shown in a similar way. \square

5.3. Conditions in terms of perfectness

A sufficient condition for our goal can be deduced from the following theorem:

Theorem 2. Let $G = (N, E), G' = (N, E')$ be two DAGs. If G' contains a subgraph H such that H is perfect and its undirected version H^\sim contains the moral graph G^M , then $\mathcal{P}^{G'} \supset \mathcal{P}^G$.

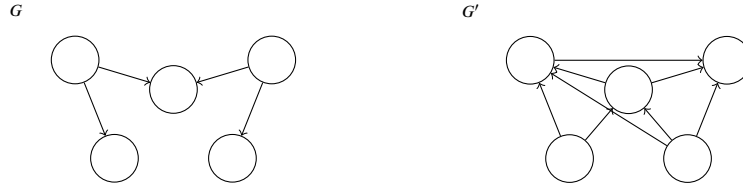


Fig. 7. Pair of DAGs G, G' that satisfies Goal II but G' does not satisfy the condition in Theorem 2. One can check that $\mathcal{P}^{G'} \supset \mathcal{P}^G$ by verifying that Condition (iv) of Theorem 1 is satisfied for all nodes.

Proof. Let $P \in \mathcal{P}^G$. By Lemma 5.9 from Cowell et al. [3] we know that P factorises undirectedly³ over the undirected graph G^M and thus over any undirected graph $L = (N, E_L)$ containing G^M . From Proposition 5.15 in Cowell et al. [3] we know that P factorises (directedly) over any perfect directed graph H such that $H^\sim = L$. Thus, $P \in \mathcal{P}^H$, and by Lemma 1, $P \in \mathcal{P}^{G'}$. Therefore when $H^\sim \supset G^M$ we have $\mathcal{P}^{G'} \supset \mathcal{P}^H \supset \mathcal{P}^G$. \square

From this theorem we can conclude that if we reverse all the edges of G , fix a topological ordering such that $\text{Leaves}(G)$ are oldest, and then add edges consonant with this topological ordering until both G' is perfect and $G' \sim \supset G^M$, we obtain an inverse of G that satisfies our goal. The example in Fig. 7 shows however that the condition that G' needs to contain a perfect subgraph H such that $H^\sim \supset G^M$ is not a necessary condition.

We do have the following necessary conditions on the graph G' to satisfy our goal:

Theorem 3. Let G, G' be DAGs. If G' is such that $\mathcal{P}^{G'} \supset \mathcal{P}^G$ and $G' \supset G^*$, then $G' \sim \supset G^M$ and for every s in N , the vertex-induced subgraph $G'[\{s\} \cup \text{des}_G(s)]$ contains a perfect subgraph H_s , such that $H_s^\sim \supset G^M[\{s\} \cup \text{des}_G(s)]$.

This theorem is based on the following two propositions. We will first prove both propositions and then show how Theorem 3 can be obtained from it.

Proposition 2. Let $\mathcal{P}^{G'} \supset \mathcal{P}^G$ and $G' \supset G^*$. Then $G' \sim \supset G^M$.

Proof. Let $t_1, t_2 \in \text{pa}_G(s)$ such that t_1 and t_2 are not joined in G and assume WLOG $\mathcal{O}'(t_2) < \mathcal{O}'(t_1)$. Assume for a contradiction that $t_2 \notin \text{pa}_{G'}(t_1)$. By Condition (iv) of Theorem 1 we must have that $t_1 \perp_G t_2 \mid \text{pa}_{G'}(t_1)$. However, since $G' \supset G^*$, we know $s \in \text{pa}_{G'}(t_1)$. Furthermore, s is a v-structure on the trail t_1, s, t_2 in G . Therefore this trail is unblocked by $\text{pa}_{G'}(t_1)$ and we arrive at a contradiction. Thus, necessarily $t_2 \rightarrow t_1$ in G' . \square

Proposition 3. If $\mathcal{P}^{G'} \supset \mathcal{P}^G$, $G' \supset G^*$, and $|\text{Roots}(G)| = 1$, then G' contains a perfect subgraph H such that $H^\sim \supset G^M$.

Proof. Our approach will be to construct a subgraph H of G' , by starting from G^* and iteratively adding edges. We show that every edge in H is also in G' and that the end result is both perfect and such that $H^\sim \supset G^M$.

Construction of H

Let \mathcal{O}' be a topological ordering of G' . We use an iterative process $H_i = (N, E_i)$ such that $H_{|N|} = H$. We start from $H_{-1} = G^*$. Then we join all parents of the same node in G according to \mathcal{O}' to obtain H_0 . In every subsequent step we join the parents in H_i of the $(i+1)$ th youngest node according to \mathcal{O}' . Formally this procedure can be defined as follows:

$$E_{-1} = E^* \tag{19}$$

$$E_0 = E_{-1} \cup E_{\text{pa}_G(G)} \tag{20}$$

$$E_{i+1} = E_i \cup E_{\text{pa}_{H_i}(r_i)}, \quad \text{for } 0 \leq i \leq |N| - 1, \tag{21}$$

with,

$$E_{\text{pa}_G(G)} = \{(t_1, t_2) : \exists s \in N \text{ such that } t_1, t_2 \in \text{pa}_G(s) \text{ and } \mathcal{O}'(t_1) < \mathcal{O}'(t_2)\} \tag{22}$$

$$E_{\text{pa}_{H_i}(r_i)} = \{(t_1, t_2) : t_1, t_2 \in \text{pa}_{H_i}(r_i), \mathcal{O}'(t_1) < \mathcal{O}'(t_2)\} \tag{23}$$

$$r_i \in N \text{ such that } \mathcal{O}'(r_i) = |N| - i. \tag{24}$$

See Fig. 8 for an example.

³ For a definition of this type of factorisation see Section 5.2 of Cowell et al. [3].

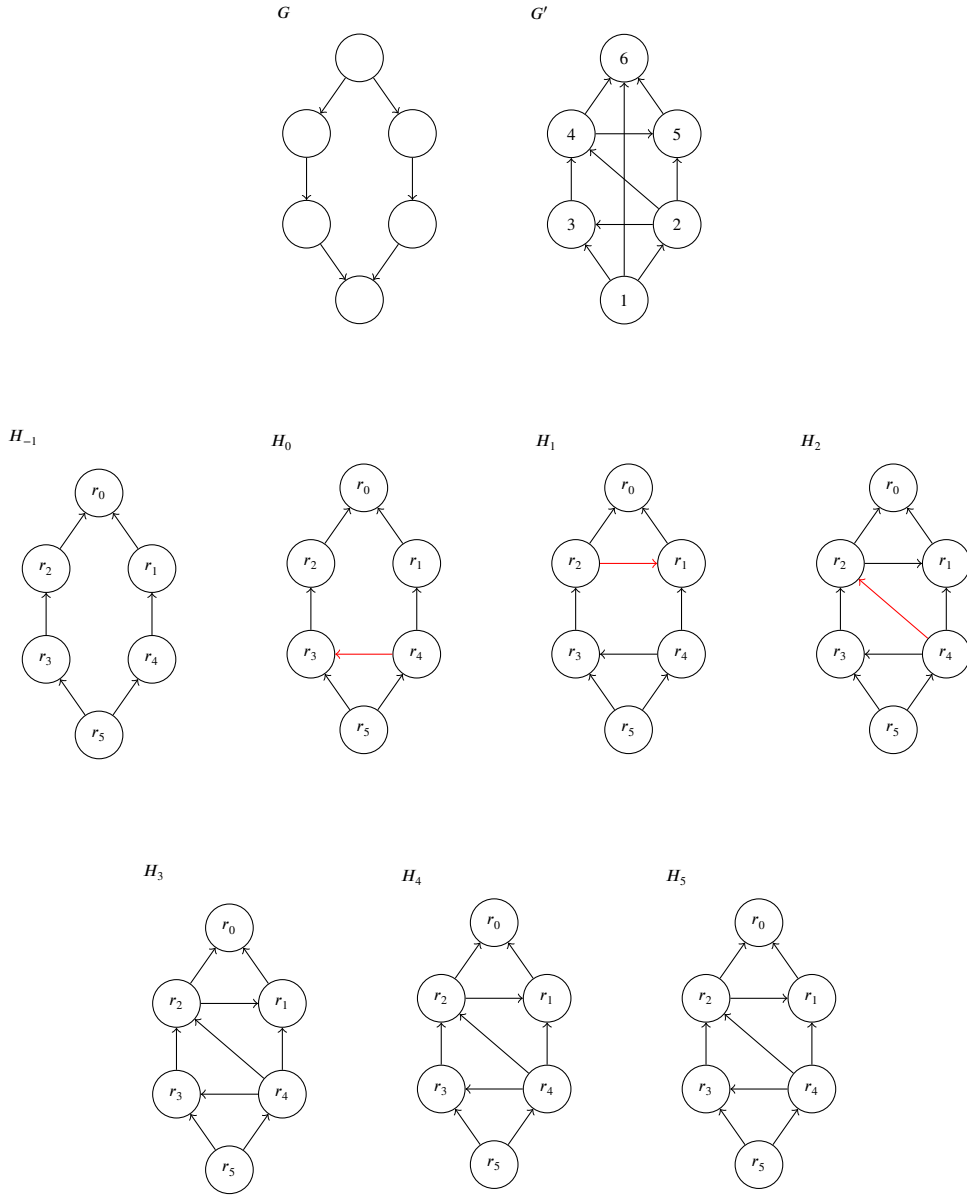


Fig. 8. (1st row) Example pair of DAGs G, G' satisfying the requirements of Proposition 3. The fact that $\mathcal{P}^{G'} \supset \mathcal{P}^G$ can be checked by verifying that Condition (iv) of Theorem 1 is satisfied for all nodes. (2nd and 3rd row) Example course of the inversion procedure (19)–(24) for the pair G, G' depicted in top row. H_{-1} is the version with the edges of G reversed. In H_0 edges $E_{\text{pa}_G(G)}$ for joining the parents in G are added. For $H_{i+1}, i \geq 0$ the parents in H_i of r_i are joined, by adding $E_{\text{pa}_{H_i}(r_i)}$ to the edge set. Note that because the parent set of r_i in H_i for $i \geq 2$ is already complete, no new edges are added in H_3, H_4 and H_5 .

H_{-1} is a subgraph of G'

This follows directly from the fact that $G' \supset G^*$.

H_0 is a subgraph of G'

This follows from Proposition 2 and the fact that both H_0 and G' are consonant with \mathcal{O}' .

H_{i+1} is a subgraph of G' for $0 \leq i \leq |N| - 1$

Lemma 4. Let $\mathcal{P}^{G'} \supset \mathcal{P}^G$, $G' \supset G^*$, $|\text{Roots}(G)| = 1$ and \mathcal{O}' a topological ordering for G' . Now let $r, u, v \in N$ such that $\mathcal{O}'(v) < \mathcal{O}'(u) < \mathcal{O}'(r)$ and $r \rightarrow v$ in G . Then, $v \rightarrow u$ in G' .

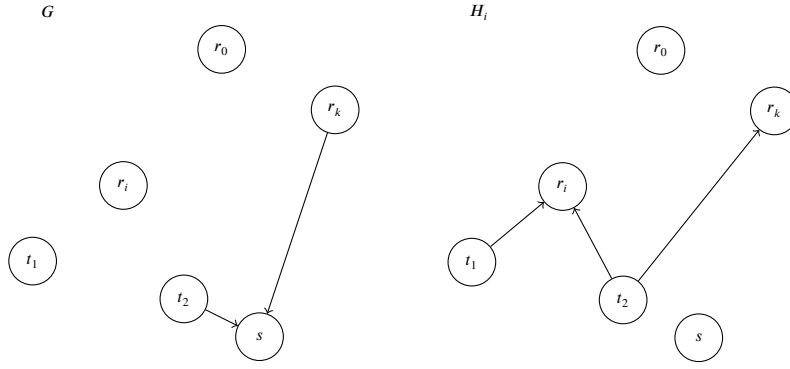


Fig. 9. Illustration of Case 2 in the proof of Proposition 3.

Proof. Let r_0 be the youngest node in \mathcal{O}' and γ_1 be a path from r_0 to u in G . Since $\mathcal{O}'(u) > \mathcal{O}'(v)$ we know $v \notin \gamma_1$. Now let γ_2 be a path in G from r_0 to r . Since $\mathcal{O}'(u) < \mathcal{O}'(r)$ we know $u \notin \gamma_2$. Now let $\gamma = \gamma_2^*; (\gamma_1, v)$ be a trail in G from u to v . Since there are no v-structures on this trail and all nodes except v are younger in \mathcal{O}' than u it follows from property (iv) of Theorem 1 that v must be a parent of u in G' . \square

Suppose $(t_2, t_1) \in E_{\text{pa}_{H_i}(r_i)}$, i.e. $t_2, t_1 \in \text{pa}_{H_i}(r_i)$ and $\mathcal{O}'(t_2) < \mathcal{O}'(t_1)$. We need to show that $t_2 \rightarrow t_1$ in G' . We distinguish two cases:

Case 1: There exists a $0 \leq j \leq i$ such that $r_j \rightarrow t_2$ in G .

This follows straight from Lemma 4 by setting $r = r_j, u = t_1, v = t_2$.

Case 2: For all $0 \leq j \leq i$ we have $r_j \not\rightarrow t_2$ in G .

See Fig. 9 for an illustration. Let $k = \min\{0 \leq j \leq i : (t_2, r_j) \in E_j\}$. Note that $(r_k, t_2) \in E_0$ follows from the definition of k by contradiction. The definition of H_0 implies the existence of s with $t_2 \rightarrow s \leftarrow r_k$ in G . The application of Lemma 4 with $r := r_k, u = t_1$ and $v = s$ gives $s \in \text{pa}_{G'}(t_1)$. Assume $t_2 \not\rightarrow t_1$ in G' for a contradiction.

We let γ_1 be a path in G from r_0 to t_1 , γ_2 a path in G from r_0 to r_k in G , and $\gamma = \gamma_2^*; \gamma_1$ the concatenation of the reversed γ_2 and γ_1 . Note that the trail (γ, s, t_2) is not blocked by $\text{pa}_{G'}(t_1)$, since for the only v-structure (r_k, s, t_2) we have $s \in \text{pa}_{G'}(t_1)$, and all other nodes on the path, except for t_2 , are younger than t_1 in \mathcal{O}' . However, by Condition (iv) of Theorem 1 we know that this trail must be blocked and we arrive at a contradiction. Therefore we conclude $t_2 \rightarrow t_1$ in G' .

H is perfect

By adding the set $E_{\text{pa}_{H_i}(r_i)}$ to the edge set, we join all the parents of r_i . After this step, no new parents of r_i are created. We perform this step for every $i \in \{0, \dots, |N| - 1\}$. Since $\{r_i : i = 0, \dots, |N| - 1\} = N$, the end result is perfect.

H^\sim contains G^M

This is ensured by adding the set $E_{\text{pa}_G(G)}$. \square

Remark 1. Note that from the proof of Proposition 3 it follows that the steps in the procedure described in (19)–(24) are actually necessary for obtaining a graph L that satisfies the following conditions:

- (i) $\mathcal{P}^L \supset \mathcal{P}^G$,
- (ii) $L \supset G^*$,
- (iii) \mathcal{O}' is a topological ordering for L .

This means that the end result H of the procedure will be a subset of any other graph L satisfying conditions (i)–(iii). By Theorem 2, H is also sufficient for satisfying these conditions. Furthermore, any other minimal inversion of G can be obtained by using the same procedure (19)–(24) but instead of \mathcal{O}' fixing a different topological ordering that is compatible with G^* .

Remark 2. Since any perfect graph with a single leave has a unique topological ordering,⁴ it follows from the proposition that any DAG G' for which there exists a DAG G such that $\mathcal{P}^{G'} \supset \mathcal{P}^G$, $G' \supset G^*$ and $\text{Roots}(G) = 1$, has a unique topological ordering as well.

Lemma 5. Let G, H be DAGs and $A \subset N$. If G, H are such that $\mathcal{P}^G \subset \mathcal{P}^H$, then the same holds for the vertex-induced subgraph of both graphs: $\mathcal{P}^{G[A]} \subset \mathcal{P}^{H[A]}$.

⁴ This follows from the fact that every node has a unique youngest parent, and the single leave is the unique youngest node of the graph.

Proof. Let \mathcal{O}_H be a topological ordering for H and $\mathcal{O}_{H[A]}$ a topological ordering for $H[A]$ that preserves the order of \mathcal{O}_H and assume that G, H are such that $\mathcal{P}^G \subset \mathcal{P}^H$. By Condition (iv) of Theorem 1 we need to show that for all $s \in A$ we have $s \perp_{G[A]} \text{pr}^{\mathcal{O}_{H[A]}(s)} | \text{pa}_{H[A]}(s)$. Therefore let $s \in A$ and $t \in \text{pr}^{\mathcal{O}_{H[A]}(s)}$ for which there exists a trail γ in $G[A]$ from s to t . We need to show that this trail is blocked by $\text{pa}_{H[A]}(s)$. Since Condition (iv) of Theorem 1 holds for the original graphs G and H , we know that this trail must be blocked by $\text{pa}_H(s)$ in G . This implies that at least one of the following conditions must hold:

- (i) There is a v-structure v on the trail such that neither v nor its descendants in G are in $\text{pa}_H(s)$.
- (ii) There is a node u on the trail that is not a v-structure and is in $\text{pa}_H(s)$.

Since $\text{pa}_{H[A]}(s) \subset \text{pa}_H(s)$ and $\text{des}_{G[A]}(v) \subset \text{des}_G(v)$, the first condition remains valid. In case of the second condition, since all nodes on the trail are in A , u must also be in $\text{pa}_{H[A]}(s)$. Therefore, if the trail is blocked by $\text{pa}_H(s)$ it is also blocked by $\text{pa}_{H[A]}(s)$. \square

Proof of Theorem 3. The first part of the theorem follows directly from Proposition 2. Note that by Lemma 5, $\mathcal{P}^{G'} \supset \mathcal{P}^G$ implies $\mathcal{P}^{G'[\{s\} \cup \text{des}_G(s)]} \supset \mathcal{P}^{G[\{s\} \cup \text{des}_G(s)]}$, for any $s \in N$. Since s is the unique root for $G[\{s\} \cup \text{des}_G(s)]$, we know from Proposition 3 that this implies that $G'[\{s\} \cup \text{des}_G(s)]$ contains a perfect subgraph H_s , such that $H_s \supset G^M[\{s\} \cup \text{des}_G(s)]$. \square

In practice, the inverse G' is often obtained by simply inverting the edges in G . In this case we have the following necessary and sufficient condition to satisfy our goal.

Theorem 4. Let G, G' be DAGs. If $G' = G^*$, then: $\mathcal{P}^{G'} \supset \mathcal{P}^G \iff \text{pa}_G(s), \text{ch}_G(s)$ are complete for all $s \in N$.

Proof. (\Leftarrow) If $\text{pa}_G(s), \text{ch}_G(s)$ are complete for all $s \in N$ and $G' = G^*$ this implies that $G' \sim \supset G^M$ and G' is perfect. The result now follows from Theorem 2.

(\Rightarrow) The fact that the parents must be complete follows from Proposition 2. Now assume by contradiction that there exists an $s \in N$ such that $u_1, u_2 \in \text{ch}_G(s)$ are not joined. Consider the distribution $P \in \mathcal{P}^G$ such that X_{u_1} and X_{u_2} are equal to X_s and all other nodes (including X_s itself) are independently distributed. Since P must factorise over G' we can write its density p as

$$p(x) = \prod_{s \in N} k^s \left(x_s | x_{\text{pa}_{G'}(s)} \right). \quad (25)$$

Since u_1 and u_2 are both independent of their parents in G' , their kernels will be simply of the form $k^{u_1} \left(x_{u_1} \right), k^{u_2} \left(x_{u_2} \right)$ respectively. This however implies that u_1 and u_2 are themselves independent which contradicts the construction of P . \square

Corollary 2. Let G be a DAG. $\mathcal{P}^G = \mathcal{P}^{G^*}$ if and only if $\text{pa}_G(s)$ and $\text{ch}_G(s)$ are complete for all $s \in N$.

Proof. This follows from twice applying Theorem 4 for $G = G$ and $G = G^*$ respectively and noting that completeness of $\text{pa}_G(s)$ and $\text{ch}_G(s)$ implies completeness of $\text{pa}_{G^*}(s)$ and $\text{ch}_{G^*}(s)$. \square

5.4. Conditions in terms of single edge operations

This paper discusses results related to the inversion of Bayesian networks in the sense of Goal II. In the proof of Proposition 3, we suggested an algorithm for inverting G , that starts by reversing all edges of G at once and then add edges where necessary. In this section we are looking at obtaining an inverse of G by reversing the edges one by one, and potentially adding edges where necessary, based on a classical result often called Meek conjecture⁵ stated below. An edge (s, t) is called *covered* when $\text{pa}_G(t) = \text{pa}_G(s) \cup \{s\}$. Note that reversing a covered edge preserves acyclicity of the graph.

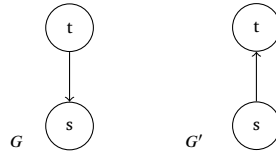
Meek [13] states the following conjecture:

Theorem 5 (Meek conjecture). Let $G = (N, E)$ and $G' = (N, E')$ be DAGs. $\mathcal{P}^{G'} \supset \mathcal{P}^G$ if and only if there exists a sequence of DAGs H_1, \dots, H_n such that $H_1 = G'$ and $H_n = G$ and H_{i+1} is obtained from H_i by one of the following operations:

- (a) covered edge reversal,
- (b) edge removal.

This result suggests the outline of an algorithm for the inversion of a Bayesian network G . This algorithm starts with $G = L_1$. L_{i+1} is obtained from L_i by performing the inverse of operation (a) or (b), i.e. going from L_{i+1} to L_i can be done by applying operation (a) or (b). Since a covered edge remains covered after reversal, this operation is equal to its inverse. The inverse of removing an edge is simply adding an edge such that acyclicity is maintained. The goal is to reverse all edges present in G without adding parents of $\text{Leaves}(G)$ that are not themselves in $\text{Leaves}(G)$.

⁵ It should be maybe more appropriately called Meek-Chickering Theorem, since it was proved in Chickering [2].

Fig. 10. Pair of graphs G, G' .

5.5. Restricting the set of possible kernel functions

The results derived in the above discuss the question what conditions G' must satisfy such that for every $P \in \mathcal{P}^G$, $\mathcal{K}^{G'}$ contains a version of the conditional distribution $P_{N \setminus \text{Leaves}(G) | \text{Leaves}(G)}$. Here it is implied that we allow for all possible kernel functions k^s in the definitions of \mathcal{P}^G and $\mathcal{K}^{G'}$. In practice, however, restrictions are often put on the space of possible kernel functions. A common choice [8] is to allow for only Gaussian kernel functions, of the form

$$k^s(x_s | x_{\text{pa}(s)}; \theta_s) \propto \exp\left(-\frac{1}{2} \left(\mu_{\theta_s}(x_{\text{pa}(s)}) - x_s\right)^2\right) \quad (26)$$

$$\mu_{\theta_s}(x_{\text{pa}(s)}) = f\left(\sum_{t \in \text{pa}(s)} \theta_{(s,t)} x_t\right), \quad (27)$$

with f some fixed possibly nonlinear function, and parametrised by the vectors $\theta_s = \{\theta_{(s,t)} \in \mathbb{R} : t \in \text{pa}(s)\}$. We will now investigate which results remain valid for the restricted case. Given a subset R of kernel functions, we will denote the restricted spaces of probability distributions and Markov kernels factorising over G by $\mathcal{P}_R^G, \mathcal{K}_R^G$ respectively. Before we dive into the results for general restrictions, we start by examining the case where R is the set of Gaussian kernel functions defined above. Consider the pair of graphs G, G' in Fig. 10. It is clear that this pair of graphs satisfies our original Goal I. However, when we restrict to the set Gaussian kernel functions, we are no longer able to model the posterior distribution exactly, as we will show now. Consider the distribution in \mathcal{P}_R^G given by

$$X_t \sim \mathcal{N}(0, 1) \quad (28)$$

$$X_s \sim \mathcal{N}(f(X_t), 1). \quad (29)$$

If the distribution $P_{t|s}$ would be in $\mathcal{K}_R^{G'}$, we would need that the joint density of X_t, X_s satisfies the following proportionality as a function of x_t

$$p(x_t, x_s) \propto \exp(-ax_t^2 + bx_t), \quad (30)$$

where only b may depend on x_s . Working out the actual joint density gives

$$p(x_t, x_s) \propto \exp\left(-\frac{1}{2}(x_t^2 + f(x_t)^2 + x_s f(x_t) + x_s^2)\right). \quad (31)$$

We can conclude that we only have that $P_{t|s} \in \mathcal{K}_R^{G'}$ if f is a linear function.⁶ From this example, we can conclude that the conditions that were sufficient for the unrestricted case, are in general not sufficient in the restricted case.

Now we look at the validity of the presented results under some general restrictions laid on the class of all kernel functions, i.e. R can be any subclass of kernel functions. We start with the equivalence of the two goals, Proposition 1. Recall that the proposition shows that finding a G' such that $\mathcal{P}^{G'} \supset \mathcal{P}^G$ and all parents in G' of nodes in $\text{Leaves}(G)$ are themselves in $\text{Leaves}(G)$, is both a necessary (statement (a)) and sufficient condition (statement (b)) to satisfy Goal I. It is easy to see that statement (b) is still valid for the restricted case, i.e. it is still a sufficient condition. However for (a) we used that when all the nodes in $\text{Roots}(G)$ are connected, any density function can be written as $p(x_{\text{Roots}(G)}) = \prod_{s \in \text{Roots}(G)} k^s(x_s | x_{\text{pa}(s)})$. This is no longer the case when we restrict the space of possible kernel functions. We have that the condition is only necessary if for every $P \in \mathcal{P}_R^G$, the marginal distribution $P_{\text{Leaves}(G)}$ factorises over a complete DAG of the leaves of G . A slightly weaker necessary condition for Goal I still holds in general, namely that for every subset $S \subset N \setminus \text{Leaves}(G)$ we need that $\mathcal{P}_R^{G'[S]} \supset \mathcal{P}_R^{G[S]}$.

For Theorem 1, note that conditions (ii) – (iv) only relate to the graph structures of G and G' . Therefore these conditions will still be equivalent for the restricted case. The implication (ii) \implies (i) does not hold in general, which was exemplified by the Gaussian kernel functions above. The implication (i) \implies (ii), on the other hand, does still hold, under the extra assumption that the restriction R is such that for any graph G , for all $A, B, S \subset N$ such that $A \not\perp_G B | S$, there is a $P \in \mathcal{P}_R^G$ for which $A \not\perp B | S$. We will sketch how this assumption is satisfied for the Gaussian kernel functions described above. Let $A, B, S \subset N$ such that $A \not\perp_G B | S$. This implies that there is a trail $\gamma = (u_0, \dots, u_n)$ in G from $u_0 = a \in A$ to $u_n = b \in B$ that is unblocked by S . We assume that the descendants of the v-structures on this path do not intersect the trail. If this is the case we replace the part of the original trail between the v-structure

⁶ Note that this remains true even when we loosen the restriction and allow $\mathcal{K}_{R_f}^{G'}$ to use a different function \tilde{f} , e.g. $\tilde{f} = f^{-1}$.

and the intersecting node by the trail going via the descendants of the former v-structure and keep doing this until all these type of v-structures are removed. Now let $\theta_{(s,t)} = 1$ for all $s = u_i, t = u_j$ with $u_i, u_j \in \gamma, u_j \in \text{pa}_G(u_i), |j - i| = 1$ and all s, t subsequent descendants of a v-structure on the trail and zero otherwise. It can be shown⁷ that for this distribution $a \perp\!\!\!\perp b \mid S$ and therefore $A \perp\!\!\!\perp B \mid S$. With this extra assumption we will now show (i) \implies (ii). Suppose $A, B, S \subset N$ such that $A \perp_{G'} B \mid S$. This implies that for all $P \in \mathcal{P}_R^{G'}$, we have $A \perp\!\!\!\perp B \mid S$. Now suppose by contradiction that $A \not\perp_G B \mid S$. By the assumption, there must be a $P \in \mathcal{P}_R^G$ for which $A \not\perp\!\!\!\perp B \mid S$, which would contradict (a). Therefore $A \perp_G B \mid S$ which shows (i) \implies (ii).

Theorem 2 is only a sufficient condition which is, by the Gaussian kernel function example, not sufficient any more in the restricted case. Theorem 3 on the other hand is only a necessary condition. The proof of this theorem only uses the necessity of the conditions in Theorem 1 which we showed above are still valid in the restricted case. We conclude that therefore Theorem 3 also still holds in the restricted case.

To conclude this section we summarise the results for the restricted case. We saw that we only have a slightly weaker necessary condition for Goal I, namely that for every subset $S \subset N \setminus \text{Leaves}(G)$ we need that $\mathcal{P}_R^{G'[S]} \supset \mathcal{P}_R^{G[S]}$. Necessary conditions for this latter condition are then provided by Theorem 1 and 3, which are still valid for the restricted case.

6. Conclusion

In this paper, we introduce some necessary and some sufficient conditions for the recognition network to be able to model the exact posterior distribution of a generative Bayesian network. In case that the generative network has a single root, the necessary and sufficient conditions coincide. However, for multiple roots there is still a gap between both conditions.

6.1. Further study directions

A further direction of study could be to find a single necessary and sufficient condition for the general case. Another interesting question is the following: “What is the smallest number of edges in an inversion G' of G ?”. Using the results on single edge operations, one could try to find an algorithm that finds an optimal inversion of G . It is generally believed that the recognition network needs many edges to make exact modelling of the posterior distribution possible (Welling, personal communication, 2022). Therefore, the number of edges in the recognition network will be reduced to make it computationally efficient. In practice, this approximation does not seem to affect the quality of the inference. This phenomenon remains poorly understood, but very relevant to the computational side of machine learning.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

No data was used for the research described in the article.

Acknowledgements

The authors would like to thank the reviewers for helpful comments. JvO would like to thank Milan Studený and Martijn Oei for helpful discussions and comments. JvO and NA acknowledge the support of the Deutsche Forschungsgemeinschaft Priority Programme “The Active Self” (SPP 2134). Lastly, JvO and PvH would like to thank Floris Triest for providing a conducive working environment.

References

- [1] R. Castelo, T. Kočka, On inclusion-driven learning of Bayesian networks, *J. Mach. Learn. Res.* 4 (Sep 2003) 527–574.
- [2] D.M. Chickering, Optimal structure identification with greedy search, *J. Mach. Learn. Res.* 3 (v) (2002) 507–554.
- [3] R.G. Cowell, P. Dawid, S.L. Lauritzen, D.J. Spiegelhalter, *Probabilistic Networks and Expert Systems*, Springer-Verlag, New York, 1999.
- [4] P. Dayan, G.E. Hinton, R.M. Neal, R.S. Zemel, The Helmholtz machine, *Neural Comput.* 7 (5) (1995) 889–904.
- [5] R.M. Dudley, *Real Analysis and Probability*, CRC Press, 2018.
- [6] I. Flesch, P.J. Lucas, Markov equivalence in Bayesian networks, in: *Advances in Probabilistic Graphical Models*, Springer, 2007, pp. 3–38.
- [7] S. Gershman, N. Goodman, Amortized inference in probabilistic reasoning, in: *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 36, 2014.
- [8] D.P. Kingma, M. Welling, Auto-encoding variational Bayes, arXiv preprint, arXiv:1312.6114, 2013.
- [9] D. Koller, N. Friedman, *Probabilistic Graphical Models: Principles and Techniques*, MIT Press, 2009.
- [10] S.L. Lauritzen, *Graphical Models*, vol. 17, Clarendon Press, 1996.

⁷ As an example, one can consider a DAG consisting of a single trail (a, \dots, b) of binary nodes for which the nodes without parents are sampled i.i.d. (Bernoulli(.5)) and the rest of the nodes are the sum of their parents modulo 2. If we condition on the nodes that are v-structures (i.e. have parents on both sides) it is easy to see that a and b are not independent.

- [11] C. Louizos, M. Welling, D.P. Kingma, Learning sparse neural networks through ℓ_0 regularization, arXiv preprint, arXiv:1712.01312, 2017.
- [12] S. Löwe, D. Madras, R. Zemel, M. Welling, Amortized causal discovery: learning to infer causal graphs from time-series data, in: Conference on Causal Learning and Reasoning, PMLR, 2022, pp. 509–525.
- [13] C. Meek, Graphical Models: Selecting causal and statistical models, PhD thesis, Carnegie Mellon University, 1997.
- [14] D. Molchanov, V. Kharitonov, A. Sobolev, D. Vetrov, Doubly semi-implicit variational inference, in: The 22nd International Conference on Artificial Intelligence and Statistics, PMLR, 2019, pp. 2593–2602.
- [15] J. Pearl, Reverend Bayes on inference engines: a distributed hierarchical approach, in: Proceedings of the Second National Conference on Artificial Intelligence, 1982, pp. 133–136.
- [16] M. Studeny, Probabilistic Conditional Independence Structures. Information Science and Statistics, Springer, London, 2005.
- [17] T. Verma, J. Pearl, Equivalence and synthesis of causal models, in: Proceedings of the Sixth Annual Conference on Uncertainty in Artificial Intelligence, 1990, pp. 255–270.
- [18] M.J. Wainwright, M.I. Jordan, et al., Graphical models, exponential families, and variational inference, Found. Trends Mach. Learn. 1 (1–2) (2008) 1–305.
- [19] S. Webb, A. Golinski, R. Zinkov, T. Rainforth, Y.W. Teh, F. Wood, et al., Faithful inversion of generative models for effective amortized inference, Adv. Neural Inf. Process. Syst. 31 (2018).