

Elizabeth S. Allman, John A. Rhodes*, Elena Stanghellini and Marco Valorta

Parameter Identifiability of Discrete Bayesian Networks with Hidden Variables

Abstract: Identifiability of parameters is an essential property for a statistical model to be useful in most settings. However, establishing parameter identifiability for Bayesian networks with hidden variables remains challenging. In the context of finite state spaces, we give algebraic arguments establishing identifiability of some special models on small directed acyclic graphs (DAGs). We also establish that, for fixed state spaces, generic identifiability of parameters depends only on the Markov equivalence class of the DAG. To illustrate the use of these results, we investigate identifiability for all binary Bayesian networks with up to five variables, one of which is hidden and parental to all observable ones. Surprisingly, some of these models have parameterizations that are generically 4-to-one, and not 2-to-one as label swapping of the hidden states would suggest. This leads to interesting conflict in interpreting causal effects.

Keywords: parameter identifiability, discrete Bayesian network, hidden variables

DOI 10.1515/jci-2014-0021

1 Introduction

A directed acyclic graph (DAG) can represent the factorization of a joint distribution of a set of random variables. To be more precise, a Bayesian network is a pair (G, P) , where G is a DAG and P is a joint probability distribution of variables in one-to-one correspondence with the nodes of G , with the property that each variable is conditionally independent of its non-descendants given its parents. It follows from this definition that the joint probability P factors according to G , as the product of the conditional probabilities of each node given its parents. Thus a discrete Bayesian network is fully specified by a DAG and a set of conditional probability tables, one for each node given its parents [1, 2].

A causal Bayesian network is a Bayesian network enhanced with a causal interpretation. Work initiated by Pearl [3, 4] investigated the identification of causal effects in causal Bayesian networks when some variables are assumed observable and others are hidden. In a non-parametric setting, with no assumptions about the state space of variables, there is a complete algorithm for determining which causal effects between variables are identifiable [5–8].

As powerful as this theory is, however, it does not address identifiability when assumptions are made on the nature of the variables. Indeed, by specializing to finite state spaces, causal effects that were non-identifiable according to the theory above may become identifiable. One particular example, with DAG shown in Figure 1, has been studied by Kuroki and Pearl [9]. If the state space of hidden variable 0 is finite, and observable variables 1 and 4 have state spaces of larger sizes, then the causal effect of variable 2 on variable 3 can be determined, for generic parameter choices.

***Corresponding author: John A. Rhodes**, Department of Mathematics and Statistics, University of Alaska Fairbanks, Fairbanks, AK, USA, E-mail: j.rhodes@alaska.edu

Elizabeth S. Allman, Department of Mathematics and Statistics, University of Alaska Fairbanks, Fairbanks, AK, USA, E-mail: esallman@alaska.edu

Elena Stanghellini, Dipartimento di Economia Finanza e Statistica, Università di Perugia, Perugia, Italy, E-mail: elena.stanghellini@stat.unipg.it

Marco Valorta, Department of Computer Science and Engineering, University of South Carolina, Columbia, SC, USA, E-mail: mgv@cse.sc.edu

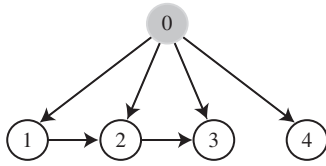


Figure 1 The DAG of a Bayesian network studied by Kuroki and Pearl [9], denoted 4-2b in the Appendix

In this paper we study in detail identification properties of certain small Bayesian networks, as a first step toward developing a systematic understanding of identification in the presence of finite hidden variables. While this includes an analysis of the model with the DAG above, our motivation is different from that of Kuroki and Pearl [9], and results were obtained independently. We make a thorough study of networks with up to five binary variables, one of which is unobservable and parental to all observable ones, as shown in Table 3 of the Appendix. These investigations lead us to develop some basic tools and arguments that can be applied more generally to questions of parameter identifiability.

In addition, for each such binary model in Table 3, we determine a value $k \in \mathbb{N} \cup \{\infty\}$ such that the marginalization from the full joint distribution to that over the observable variables is generically k -to-one. Although we restrict this exhaustive study to binary models for simplicity, straightforward modifications to our arguments would extend them to larger state spaces. A typical requirement for such an extended identifiability result is that the state spaces of observable variables be sufficiently large, relative to that of the hidden variable, as in the result of Kuroki and Pearl [9] described earlier. Interestingly, that result restricted to finite state spaces follows easily from our framework, and can be obtained for continuous state spaces of observable variables using arguments of Allman et al. [10].

We use the term “DAG model” for the collection of all Bayesian networks with the same DAG and specification of state spaces for the variables. With the conditional probability tables of nodes given their parents forming the parameters of the model, we thus allow these tables to range over all valid tables of a fixed size to give the parameter space of such a model.

That some of the DAG models we consider have non-identifiable parameters is a consequence of the well-known non-uniqueness (in most circumstances) of non-negative rank decompositions of matrices. An example is the infinite-to-one parameterization of model 4-2a in Table 3. For greater detail on this issue see the work of Mond et al. [11] and Kubjas et al. [12].

In dealing with discrete unobserved variables, another well-understood identifiability issue is sometimes called *label swapping*. If the latent variable has r states, there are $r!$ parameter choices, obtained by permuting the state labels of the latent variable, that generate the same observable distribution. Thus the parameterization map is generically at least $r!$ -to-one. For models with a single binary latent variable, it is thus commonly expected that parameterizations are either infinite-to-one due to a parameter space of too high a dimension, or 2-to-one due to label swapping. Our work, however, finds surprisingly simple examples such that the mapping is 4-to-one, so that subtler non-identifiability issues arise.

Our analysis arises from an algebraic viewpoint of the identifiability problem. With finite state spaces the parameterization maps for DAG models with hidden variables are polynomial. Given a distribution arising from the model, the parameters are identifiable precisely when a certain system of multivariate polynomial equations has exactly one solution (up to label swapping of states for hidden variables). Though in principle computational algebra software can be used to investigate parameter identifiability, the necessary calculations are usually intractable for even moderate size DAGs and/or state spaces. In addition, one runs into issues of complex versus real roots, and the difficulty of determining when real roots lie within stochastic bounds. While our arguments are fundamentally algebraic, they do not depend on any machine computations.

If a single polynomial $p(x)$ in one variable is given, of degree n , then it is well known that the map from \mathbb{C} to \mathbb{C} that it defines will be generically n -to-one. Indeed the equation $p(x) = a$ will be of degree n for each choice of a , and generically will have n distinct roots. This fact generalizes to polynomial maps from

\mathbb{C}^n to \mathbb{C}^m ; there always exists a $k \in \mathbb{N} \cup \{\infty\}$ such that the map is generically k -to-one. However if $p(x)$ has real coefficients, and is instead viewed as a map from (a subset of) \mathbb{R} to \mathbb{R} , it may not have a generic k -to-one behavior. For instance, from a typical graph of a cubic one sees there can be sets of positive measure on which it is 3-to-one, and others on which it is one-to-one, as well as an exceptional set of measure zero on which the cubic is 2-to-one. While this exceptional set arises since a polynomial may have repeated roots, the lack of a generic k -to-one behavior is due to passing from considering a complex domain for the function, to a real one.

The fact that the polynomial parameterizations for the models investigated here have a generic k -to-one behavior on their parameter space thus depends on the particular form of the parameterizations. For those binary models in Table 3, we prove this essentially one model at a time, while obtaining the value for k . In the case of finite k , our arguments actually go further and characterize the k elements of $\phi^{-1}(\phi(\theta))$ in terms of a generic θ . Of course when $k = 2$ this is nothing more than label swapping, but for the cases of $k = 4$ more is required. Precise statements appear in later sections. In some cases, we also give descriptions of an exceptional subset of Θ where the generic behavior may not hold. In all cases, the reader can deduce such a set from our arguments.

After setting terminology in Section 2, in Section 3 we establish that, when all variables have fixed finite state spaces, Markov equivalent DAGs specify parameter equivalent models. Thus in answering generic identifiability questions one need only consider Markov equivalence classes of DAGs. In Section 4 we revisit the fundamental result due to Kruskal [13], as developed in Allman et al. [10] for identifiability questions. We give explicit identifiability procedures for the DAG this to which this result applies most directly (model 3-0), and also for the DAG of model 4-3b. These two DAGs are basic cases whose known identifiability is then leveraged in Section 5 to determine generic identifiability results for all the binary DAGs we catalog. Although in Section 5 we do not push our arguments toward exhaustive consideration of non-binary models, in many cases it would be straightforward to do so. For instance, if all variables associated to a DAG have the same size state space, little in our arguments needs to be modified.

Finally, in Section 6 we construct an explicit distribution for the generically 4-to-one parameterization of model 4-3e in which there are two different causal effects consistent with the observable distribution. This is possible because the parameter sets that give rise to this distribution differ in ways beyond label swapping. Determining causal effects coherently in this context is thus impossible. This example provokes a general caution: The parameterization of a discrete DAG model can be k -to-one with k larger than one would expect from label swapping, and when this occurs quantifying causal effects can be highly problematic.

We view the main contribution of this paper not as the determination of parameter identifiability for the specific binary models we consider, but rather as the development of the techniques by which we establish our results. Ultimately, one would like fairly simple graphical rules to determine which parameters are identifiable, and perhaps even to yield formulas for them in terms of the joint distribution. Establishing similar results for more general graphical models, not specified by a DAG, is also desirable. Some work in this context already exists (see, e.g., Stanghellini and Vantaggi [14]).

2 Discrete DAG models and parameter identifiability

The models we consider are specified in part by DAGs $\mathcal{G} = (V, E)$ in which nodes $v \in V$ represent random variables X_v , and directed edges in E imply certain independence statements for the joint distribution of all variables [15]. A bipartition of $V = O \sqcup H$ is given, in which variables associated to nodes in O or H are observable or hidden, respectively. Finally, we fix finite state spaces, of size n_v for each variable X_v .

A DAG \mathcal{G} entails a collection of conditional independence statements on the variables associated to its nodes, via d-separation, or an equivalent separation criterion in terms of the moral graph on ancestral sets.

A joint distribution of variables satisfies these statements precisely when it has a factorization according to \mathcal{G} as

$$P = \prod_{v \in V} P(X_v | X_{\text{pa}(v)}),$$

with $\text{pa}(v)$ denoting the set of parents of v in \mathcal{G} . We refer to the conditional probabilities $\theta = (P(X_v | X_{\text{pa}(v)}))_{v \in V}$ as the *parameters* of the DAG model, and denote the space of all possible choices of parameters by $\Theta = \Theta_{\mathcal{G}, \{n_v\}}$. The parameterization map for the joint distribution of all variables, both observable and hidden, is denoted as

$$\phi : \Theta \rightarrow \Delta \left(\prod_{v \in V} n_v \right)^{-1},$$

where Δ^k is the k -dimensional probability simplex of stochastic vectors in \mathbb{R}^{k+1} . Thus $\phi(\Theta)$ is precisely the collection of all probability distributions satisfying the conditional independence statements associated to \mathcal{G} (and possibly additional ones).

Since the probability distribution for the model with hidden variables is obtained from that of the fully observable model, its parameterization map is

$$\phi^+ = \sigma \circ \phi : \Theta \rightarrow \Delta \left(\prod_{v \in O} n_v \right)^{-1},$$

where σ denotes the appropriate map marginalizing over hidden variables. The set $\phi^+(\Theta)$ is thus the collection of all observable distributions that arise from the hidden variable model. This collection depends not only on the DAG and designated state spaces of observable variables, but also on the state spaces of hidden variables, even though the sizes of hidden state spaces are not readily apparent from an observable joint distribution.

With all variables having finite state spaces, the parameter space Θ can be identified with the closure of an open subset of $[0, 1]^L$, for some L . We refer to L as the dimension of the parameter space. The dimension of Θ is easily seen to be

$$\dim(\Theta) = \sum_{v \in V} \left((n_v - 1) \prod_{w \in \text{pa}(v)} n_w \right). \tag{1}$$

In the case of all binary variables, this simplifies to

$$\dim(\Theta) = \sum_{v \in V} 2^{|\text{pa}(v)|} = \sum_{k=0}^{\infty} m_k 2^k, \tag{2}$$

where m_k is the number of nodes in \mathcal{G} with in-degree k .

If a statement is said to hold for *generic parameters* or *generically* then we mean it holds for all parameters in a set of the form $\Theta \setminus E$, where the exceptional set E is a proper algebraic subset of Θ . (Recall an *algebraic subset* is the zero set of a finite collection of multivariate polynomials.) As proper algebraic subsets of \mathbb{R}^n are always of Lebesgue measure zero, a statement that holds generically can fail only on a set of measure zero.

As an example of this language, for any DAG model with all variables finite and observable, generic parameters lead to a distribution faithful to the DAG, in the sense that those conditional independence statements implied by d-separation rules will hold, and no others [16]. Equivalently, a generic distribution from such a model is faithful to the DAG.

There are several notions of identifiability of parameters of a model; we refer the reader to Allman et al. [10]. The strictest notion, that the parameterization map is one-to-one, is easily seen to hold when all DAG variables are observable with mild additional assumptions (e.g., positivity of all parameters). If a model has

hidden variables, then this is too strict a notion of identifiability, as the well-known issue of label swapping arises: One can permute the names of the states of hidden variables, making appropriate changes to associated parameters, without changing the joint distribution of the observable variables. For a model with one r -state hidden variable, label swapping implies that for any generic $\theta_1 \in \Theta$ there are at least $r! - 1$ other points $\theta_j \in \Theta$ with $\phi^+(\theta_1) = \phi^+(\theta_j)$. But since these are isolated parameter points that differ only by state labeling, this issue does not generally limit the usefulness of a model, provided that we remain aware of it when interpreting parameters.

The strongest useful notion of identifiability for models with hidden variables is that for generic $\theta_1 \in \Theta$, if $\phi^+(\theta_1) = \phi^+(\theta_2)^+$, then θ_1 and θ_2 differ only up to label swapping for hidden variables. This notion is our primary focus in this paper, which we refer to as *generic identifiability up to label swapping*. In particular, for models with a single binary hidden variable it is equivalent to the parameterization map being generically 2-to-one.

3 Markov equivalence and parameter identifiability

Two DAGs on the same sets of observable and hidden nodes are said to be *Markov equivalent* if they entail the same conditional independence statements through d-separation. (Note this notion does not distinguish between observable and hidden variables; all are treated as observable.) Thus for fixed choices of state spaces of the variables, two different but Markov equivalent DAGs, $\mathcal{G}_1 \cong \mathcal{G}_2$, have different parameter spaces Θ_1, Θ_2 , and different parameterization maps, yet $\phi_1(\Theta_1) = \phi_2(\Theta_2)$.

For studying identifiability questions, it is helpful to first explore the relationship between parameterizations for Markov equivalent graphs. A simple example, with no hidden variables, is instructive. Consider the DAGs on two observable nodes

$$1 \rightarrow 2, \quad 1 \leftarrow 2,$$

which are equivalent, since neither entails any independence statements. Now the particular probability distribution $P(X_1 = i, X_2 = j) = P_{ij}$ with

$$P = \begin{pmatrix} 1/2 & 0 \\ 1/2 & 0 \end{pmatrix}$$

requires parameters on the first DAG to be

$$P(X_1) = (1/2, 1/2), \quad P(X_2|X_1) = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix},$$

while parameters on the second DAG can be

$$P(X_2) = (1, 0), \quad P(X_1|X_2) = \begin{pmatrix} 1/2 & 1/2 \\ t & 1-t \end{pmatrix}$$

for any $t \in [0, 1]$. Thus this particular distribution has identifiable parameters for only one of these DAGs. (Here and in the rest of the paper conditional probability tables specifying parameters have rows corresponding to states of conditioning, i.e., parent, variables.)

Of course, this probability distribution was a special one, and is atypical for these models, which are easily seen to have generically identifiable parameters (as do all DAG models without hidden variables). Nonetheless, it illustrates the need for “generic” language and careful arguments for results such as the following.

Theorem 1. *With all variables having fixed finite state spaces, consider two Markov equivalent DAGs, \mathcal{G}_1 and \mathcal{G}_2 , possibly with hidden nodes. If the parameterization map ϕ_1^+ is generically k -to-one for some $k \in \mathbb{N}$, then ϕ_2^+ is also generically k -to-one.*

In particular if such a model has parameters that are generically identifiable up to label swapping, so does every Markov equivalent model.

This theorem is a consequence of the following:

Lemma 2. *With all variables having finite state spaces, consider two Markov equivalent DAGs, \mathcal{G}_1 and \mathcal{G}_2 , with parameter spaces Θ_i and parameterization maps ϕ_i , $i \in \{1, 2\}$, for the joint distribution of all variables. Then there are generic subsets $S_i \subseteq \Theta_i$ and a rational homeomorphism $\psi : S_1 \rightarrow S_2$, with rational inverse, such that for all $\theta \in S_1$*

$$\phi_1(\theta) = \phi_2(\psi(\theta)).$$

Proof. Recall that an edge $i \rightarrow j$ of a DAG is said to be covered if $\text{pa}(j) = \text{pa}(i) \cup \{i\}$. By Chickering [17], Markov equivalent DAGs differ by applying a sequence of reversals of covered edges.

We thus first assume the \mathcal{G}_i differ by the reversal of a single covered edge $i \rightarrow j$ of \mathcal{G}_1 . Let $W = \text{pa}_{\mathcal{G}_1}(i) = \text{pa}_{\mathcal{G}_2}(j)$, so $\text{pa}_{\mathcal{G}_1}(j) = W \cup \{i\}$, $\text{pa}_{\mathcal{G}_2}(i) = W \cup \{j\}$. Now any $\theta \in \Theta_1$ is a collection of conditional probabilities $P(X_v|X_{\text{pa}(v)})$, including $P(X_i|W)$, $P(X_j|X_i, W)$. From these, successively define

$$P(X_i, X_j|W) = P(X_j|X_i, W)P(X_i|W),$$

$$P(X_j|W) = \sum_k P(X_i = k, X_j|W),$$

$$P(X_i|X_j, W) = P(X_i, X_j|W)/P(X_j|W).$$

Using these last two conditional probabilities, along with those specified by θ for all $v \neq i, j$, define parameters $\psi(\theta) \in \Theta_2$. Now ψ is defined and continuous on the set S_1 where $P(X_i|W)$ and $P(X_j|X_i, W)$ are strictly positive.

One easily checks that the same construction applied to the edge $j \rightarrow i$ in \mathcal{G}_2 gives the inverse map.

If $\mathcal{G}_1, \mathcal{G}_2$ differ by a sequence of edge reversals, one defines the S_i as subsets where all parameters related to the reversed edges are strictly positive, and let ψ be the composition of the maps for the individual reversals. \square

Proof of Theorem 1. Suppose Θ_1 has a generic subset S on which ϕ_1^+ is k -to-one and the map ψ of Lemma 2 is invertible. Then $\psi(S)$ will be a generic subset of Θ_2 , and the identity

$$\phi_2^+(\theta) = \phi_1^+(\psi^{-1}(\theta))$$

from Lemma 2 shows that ϕ_2^+ is k -to-one on $\psi(S)$. Thus we need only establish the existence of such an S .

Let $S_1 = \Theta_1 \setminus E_1$, $S_2 = \Theta_2 \setminus E_2$ be the generic sets of Lemma 2. Let $S'_1 = \Theta_1 \setminus E'_1$ be a generic set on which ϕ_1^+ is k -to-one. We may thus assume E_1, E'_1, E_2 are all proper algebraic subsets. Since ϕ_1^+ is generically k -to-one with finite k , the set $(\phi_1^+)^{-1}(\phi_1^+(E_1))$ must be contained in a proper algebraic subset of Θ_1 , say E''_1 . We may therefore take $S = \Theta_1 \setminus (E'_1 \cup E''_1)$. \square

4 Two special models

In this section, we explain how one may explicitly solve for parameter values from a joint distribution of the observable variables for models specified by two specific DAGs with hidden nodes.

Parameter identifiability of the model with DAG shown in Figure 2 is an instance of a more general theorem of Kruskal [13] (see also [18, 19]). However, known proofs of the full Kruskal theorem do not yield an explicit procedure for recovering parameters. Nonetheless, a proof of a restricted theorem (the essential idea of which is not original to this work, and has been rediscovered several times) does. We include this argument for Theorem 3, since it is still not widely known and provides motivation for the approach to the proof of Theorem 4 for models associated to a second DAG, shown in Figure 3. Our analysis of the second model appears to be entirely novel. For both models, we characterize the exceptional parameters for which these procedures fail, giving a precise characterization of a set containing all non-identifiable parameters.

4.1 Explicit cases of Kruskal's theorem

The model we consider has the DAG of model 3-0 in Table 3, also shown in Figure 2 for convenience.

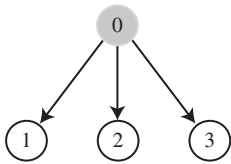


Figure 2 The DAG of model 3-0, the Kruskal model

Parameters for the model are:

1. $\mathbf{p}_0 = P(X_0) \in \Delta^{n_0-1}$, a stochastic vector giving the distribution for the n_0 -state hidden variable X_0 .
2. For each of $i = 1, 2, 3$, a $n_0 \times n_i$ stochastic matrix $M_i = P(X_i|X_0)$.

We use the following terminology.

Definition. The *Kruskal row rank* of a matrix M is the maximal number r such that every set of r rows of M is linearly independent.

Note that the Kruskal row rank of a matrix may be less than its rank, which is the maximal r such that *some* set of r rows is independent.

Our special case of Kruskal's theorem is the following:

Theorem 3. Consider the model represented by the DAG of model 3-0, where variables X_i have $n_i \geq 2$ states, with $n_1, n_2 \geq n_0$. Then generic parameters of the model are identifiable up to label swapping, and an algebraic procedure for determination of the parameters from the joint probability distribution $P(X_1, X_2, X_3)$ can be given.

More specifically, if \mathbf{p}_0 has no zero entries, M_1, M_2 have rank n_0 , and M_3 has Kruskal row rank at least 2, then the parameters can be found through determination of the roots of certain n_0 th degree univariate polynomials and solving linear equations. The coefficients of these polynomials and linear systems are rational expressions in the joint distribution.

Proof. For simplicity, consider first the case $n_0 = n_1 = n_2 = n$. Let $P = P(X_1, X_2, X_3)$ be a probability distribution of observable variables arising from the model, viewed as a $n \times n \times n_3$ array.

Marginalizing P over X_3 (i.e., summing over the 3rd index), we obtain a matrix which, in terms of the unknown parameters, is the matrix product

$$P_{..+} = P(X_1, X_2) = M_1^T \text{diag}(\mathbf{p}_0) M_2.$$

Similarly, if $M_3 = (m_{ij})$, then the slices of P with third index fixed at i (i.e., the conditional distributions given $X_i = i$, up to normalization) are

$$P_{\cdot,i} = P(X_1, X_2, X_3 = i) = M_1^T \text{diag}(\mathbf{p}_0) \text{diag}(M_3(\cdot, i)) M_2,$$

where $M_3(\cdot, i)$ is the i th column of M_3 .

Assuming M_1, M_2 are non-singular, and \mathbf{p}_0 has no zero entries, $P_{\cdot,+}$ is invertible and we see

$$P_{\cdot,+}^{-1} P_{\cdot,i} = M_2^{-1} \text{diag}(M_3(\cdot, i)) M_2. \tag{3}$$

Thus the entries of the columns of M_3 can be determined (without order) by finding the eigenvalues of the $P_{\cdot,+}^{-1} P_{\cdot,i}$, and the rows of M_2 can be found by computing the corresponding left eigenvectors, normalizing so the entries add to 1. (If M_3 has repeated entries in the i th column, the eigenvectors may not be uniquely determined. However, since the matrices $P_{\cdot,+}^{-1} P_{\cdot,i}$ for various i commute, and M_3 has Kruskal row rank 2 or more, the set of these matrices do uniquely determine a collection of simultaneous 1-dimensional eigenspaces. We leave the details to the reader.) This determines M_2 and M_3 , up to the simultaneous ordering of their rows.

A similar calculation with $P_{\cdot,i} P_{\cdot,+}^{-1}$ determines M_1 , and M_3 , up to the row order. Since the rows of M_3 are distinct (because it has Kruskal rank 2), fixing some ordering of them fixes a consistent order of the rows of all of the M_i .

Finally, one determines \mathbf{p}_0 from $M_1^{-T} P_{\cdot,+} M_2^{-1} = \text{diag}(\mathbf{p}_0)$.

The hypotheses on the rank and Kruskal rank of the parameter matrices can be expressed through the non-vanishing of minors, so all assumption on parameters used in this procedure can be phrased as the non-vanishing of certain polynomials. As a result, the exceptional set where it cannot be performed is contained in a proper algebraic subset of the parameter set.

Since the computations to perform the procedure involve computing eigenvalues and eigenvectors of matrices whose entries are rational in the joint distribution, the second paragraph of the theorem is justified.

In the more general case of $n_1, n_2 \geq n_0$, one can apply the argument above to $n_0 \times n_0 \times n_3$ subarrays of P corresponding to submatrices of M_1 and M_2 that are invertible. All such subarrays will lead to the same eigenvalues of the matrices analogous to those of eq. (3), so eigenvectors can be matched up to reconstruct entire rows of M_1 and M_2 . The vector \mathbf{p}_0 is determined by a formula similar to that above, using a subarray of the marginalization $P_{\cdot,+}$. □

4.2 Another special model

The model we consider next has the DAG of model 4-3b in Table 3, reproduced in Figure 3 for convenience.

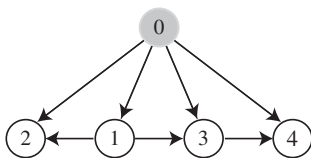


Figure 3 The DAG of model 4-3b

Parameters for the model are:

1. $\mathbf{p}_0 = P(X_0) \in \Delta^{n_0-1}$, a stochastic vector giving the distribution for the n_0 -state hidden variable X_0 .
2. Stochastic matrices $M_1 = P(X_1|X_0)$ of size $n_0 \times n_1$; $M_i = P(X_i|X_0, X_1)$ of size $n_0 n_1 \times n_i$ for $i = 2, 3$; and $M_4 = P(X_4|X_0, X_3)$ of size $n_0 n_3 \times n_4$.

Theorem 4. Consider the model represented by the DAG of model 4-3b, where variables X_i have $n_i \geq 2$ states, with $n_2, n_4 \geq n_0$. Then generic parameters of the model are identifiable up to label swapping, and an algebraic procedure for determination of the parameters from the joint probability distribution $P(X_1, X_2, X_3, X_4)$ can be given.

More specifically, suppose \mathbf{p}_0, M_1, M_3 have no zero entries, the $n_0 \times n_2$ and $n_0 \times n_4$ matrices

$$M_2^i = P(X_2|X_0, X_1 = i), \quad 1 \leq i \leq n_1, \quad \text{and}$$

$$M_4^j = P(X_4|X_0, X_3 = j), \quad 1 \leq j \leq n_3$$

have rank n_0 , and there exists some i, i' with $1 \leq i < i' \leq n_1$ such that for all $1 \leq j < j' \leq n_3$, $1 \leq k < k' \leq n_0$ the entries of M_3 satisfy inequality (7). Then from the resulting joint distribution the parameters can be found through determination of the roots of certain n th degree univariate polynomials and solving linear equations. The coefficients of these polynomials and linear systems are rational expressions in the entries of the joint distribution.

Proof. Consider first the case $n_0 = n_2 = n_4 = n$. With $P = P(X_1, X_2, X_3, X_4)$ viewed as an $n_1 \times n \times n_3 \times n$ array, we work with $n \times n$ slices of P ,

$$P_{ij} = P(X_1 = i, X_2, X_3 = j, X_4),$$

that is, we essentially condition on X_1, X_3 , though omit the normalization.

Note that these slices can be expressed as

$$P_{ij} = (M_2^i)^T D_{ij} M_4^j, \quad (4)$$

where $D_{ij} = \text{diag}(P(X_0, X_1 = i, X_3 = j))$ is the diagonal matrix given in terms of parameters by

$$D_{ij}(k, k) = \mathbf{p}_0(k) M_1(k, i) M_3((k, i), j),$$

and M_2^i and M_4^j are as in the statement of the theorem.

Equation (4) implies for $1 \leq i, i' \leq n_1$ and $1 \leq j, j' \leq n_3$ that

$$P_{ij}^{-1} P_{i'j'} P_{i'j}^{-1} P_{ij} = (M_4^j)^{-1} D_{ij}^{-1} D_{i'j'} D_{ij}^{-1} D_{i'j'} M_4^j, \quad (5)$$

and the hypotheses on the parameters imply the needed invertibility. But this shows the rows of M_4^j are left eigenvectors of this product.

In fact, if $i \neq i', j \neq j'$, then the eigenvalues of this product are distinct, for generic parameters. To see this, note the eigenvalues are

$$M_3((k, i), j') M_3((k, i'), j) / (M_3((k, i), j) M_3((k, i'), j')), \quad (6)$$

for $1 \leq k \leq n$, so distinctness of eigenvalues is equivalent to

$$\begin{aligned} & M_3((k, i), j') M_3((k, i'), j) M_3((k', i), j) M_3((k', i'), j') \\ & \neq M_3((k, i), j) M_3((k, i'), j') M_3((k', i), j') M_3((k', i'), j), \end{aligned} \quad (7)$$

for all $1 \leq k < k' \leq n$. Thus a generic choice of M_3 leads to distinct eigenvalues.

With distinct eigenvalues, the eigenvectors are determined up to scaling. But since each row of M_4^j must sum to 1, the rows of M_4^j are therefore determined by P .

The ordering of the rows of the M_4^j has not yet been determined. To do this, first fix an arbitrary ordering of the rows of M_4^1 , say, which imposes an arbitrary labeling of the states for X_0 . Then using eq. (4), from $P_{i,1} (M_4^1)^{-1}$ we can determine $D_{i,1}$ and M_2^i with their rows ordered consistently with M_4^1 . For $j \geq 1$, using eq. (4) again, from $(M_2^i)^{-T} P_{ij}$ we can determine D_{ij} and M_4^j with a consistent row order. Thus M_2 and M_4 are determined.

To determine the remaining parameters, again appealing to eq. (4), we can recover the distribution $P(X_0, X_1, X_2)$ using

$$(M_2^i)^{-T} P_{i,j} (M_4^j)^{-1} = \text{diag}(P(X_0, X_1 = i, X_3 = j)).$$

With X_0 no longer hidden, it is straightforward to determine the remaining parameters.

The general case of $n_0 \leq n_2, n_4$ is handled by considering subarrays, just as in the proof of the preceding theorem. \square

Remark. In the case of all binary variables, the expression in eq. (6) is just the conditional odds ratio for the observed variables X_1, X_3 , conditioned on X_0 . Inequality (7) can thus be interpreted as saying there is a non-zero 3-way interaction between the variables X_0, X_1, X_2 , which is the generic situation.

5 Small binary DAG models

All variables are assumed binary throughout this section. In Table 3 of the Appendix, we list each of the binary DAG models with one latent node which is parental to up to 4 observable nodes. We number the graphs as A - Bx where $A = |O| = |V| - 1$ is the number of observed variables, $B = |E| - |O|$ is the number of directed edges between the observed variables, and x is a letter appended to distinguish between several graphs with these same features. As the table presents only the case that all variables are binary, the observable distribution lies in a space of dimension $2^A - 1$.

The primary information in this table is in the column for k , indicating the parameterization map is generically k -to-one. As discussed in the introduction, the existence of such a k is not obvious, and does not follow from the behavior of general polynomial maps in real variables.

The models 4-3e and 4-3f, for which the parameterization maps are generically 4-to-one, are particularly interesting cases, as for these models there are non-identifiability issues that arise neither from overparameterization (in the sense of a parameter space of larger dimension than the distribution space) nor from label swapping. While these models are ones that can plausibly be imagined as being used for data analysis, they have a rather surprising failure of identifiability, which is explored more precisely in Section 6.

We now turn to establishing the results in Table 3.

For many of the models A - Bx the dimension of the parameter space computed by eq. (2) exceeds the dimension $2^A - 1$ of the probability simplex in which the joint distribution of observed variables lies. In these cases, the following proposition applies to show the parameterization is generically infinite-to-one. We omit its proof for brevity.

Proposition 5. *Let $f : S \rightarrow \mathbb{R}^m$ be any map defined by real polynomials, where S is an open subset of \mathbb{R}^n and $n > m$. Then f is generically infinite-to-one.*

This proposition applies to all models in Table 3 with an infinite-to-one parameterization, with the single exception of 4-2a. For that model, amalgamating X_1 and X_2 together, and likewise X_3 and X_4 , we obtain a model with two 4-state observed variables that are conditionally independent given a binary hidden variable X_0 . One can show that the probability distributions for this model form an 11-dimensional object, and then a variant of Proposition 5 applies.

For models 3-0 and 4-3b (and the Markov equivalent 4-3a), specializing Theorems 3 and 4 of the previous section to binary variables yields the claims in the table.

For the remaining models, the strategy is to first marginalize or condition on an observable variable to reduce the model to one already understood. One then attempts to “lift” results on the reduced model back to the original one.

We consider in detail only some of the models, indicating how the arguments we give can be adapted to others with minor modifications.

5.1 Model 4-1

Referring to Figure 4, since node 2 is a sink, marginalizing over X_2 gives an instance of model 3-0 with the same parameters, after discarding $P(X_2|X_0, X_1)$. Thus generically all parameters except $P(X_2|X_0, X_1)$ are determined, up to label swapping.

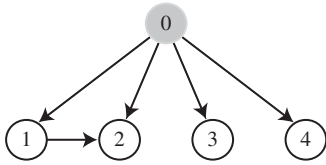


Figure 4 The DAG of model 4-1

But note that if the (unknown) joint distribution of X_0, X_1, X_2, X_3 is written as an 8×2 matrix U , with

$$U((i, j, k), \ell) = P(X_0 = \ell, X_1 = i, X_2 = j, X_3 = k),$$

and $M_4 = P(X_4|X_0)$, then the matrix product UM_4 has entries

$$(UM_4)((i, j, k), \ell) = P(X_1 = i, X_2 = j, X_3 = k, X_4 = \ell),$$

which form the observable joint distribution. Since generically M_4 is invertible, from the observable distribution and each of the already identified label swapping variants of M_4 we can find U . From U we marginalize to obtain $P(X_0, X_1, X_2)$ and $P(X_0, X_1)$. Under the generic condition that $P(X_0), P(X_1|X_0)$ are strictly positive, $P(X_0, X_1)$ is as well, and so we can compute $P(X_2|X_0, X_1) = P(X_0, X_1, X_2)/P(X_0, X_1)$.

Models 4-0 and 4-2d are handled similarly, by marginalizing over the sink nodes 4 and 3, respectively.

An alternative argument for models 4-1 and 4-0 proceeds by amalgamating the observed variables, X_1, X_2 , into a single 4-state variable, and applying Theorem 3 directly to that model. We leave the details to the reader.

5.2 Models 4-2b,c

Up to renaming of nodes, the DAGs for models 4-2b and 4-2c are Markov equivalent. Thus by Theorem 1, it is enough to consider model 4-2c, as shown in Figure 5.

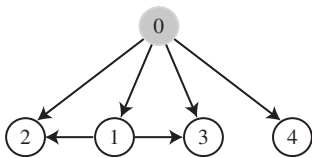


Figure 5 The DAG of model 4-2c

We condition on $X_1 = j, j = 1, 2$ to obtain two related models. Letting $X_i^{(j)}$ denote the conditioned variable at node i , the resulting observable distributions are

$$\begin{aligned} P(X_2^{(j)}, X_3^{(j)}, X_4^{(j)}) &= P(X_2, X_3, X_4|X_1 = j) \\ &= P(X_1 = j)^{-1}P(X_1 = j, X_2, X_3, X_4). \end{aligned}$$

With a hidden variable $X_0^{(j)}$ and observed variables $X_2^{(j)}, X_3^{(j)}, X_4^{(j)}$, these distributions arise from a DAG like that of model 3-0. With parameters for the original model $\mathbf{p}_0 = P(X_0)$, 2×2 matrices $M_i = P(X_i|X_0)$ for $i = 1, 4$, and 4×2 matrices $M_i = P(X_i|X_0, X_1)$, $i = 2, 3$ and \mathbf{e}_j the standard basis vector, parameters for the conditioned models are

1. the vector

$$\begin{aligned} \mathbf{p}_0^{(j)} &= P(X_0^{(j)}) = P(X_0|X_1 = j) \\ &= P(X_1 = j)^{-1}P(X_0, X_1 = j) \\ &= \frac{1}{\mathbf{p}_0^T M_1 \mathbf{e}_j} (\text{diag}(\mathbf{p}_0) M_1 \mathbf{e}_j), \end{aligned}$$

2. the 2×2 stochastic matrix $M_4^{(i)} = P(X_4^{(i)}|X_0^{(i)}) = M_4$, and
 3. for $i = 2, 3$, the 2×2 stochastic matrix $M_i^{(j)} = P(X_i^{(j)}|X_0^{(j)})$, whose rows are the $(0, j)$ and $(1, j)$ rows of M_i .

Now if \mathbf{p}_0 and column j of M_1 have non-zero entries, it follows that $\mathbf{p}_0^{(j)}$ has no zero entries. If additionally $M_2^{(j)}, M_3^{(j)}, M_4$ all have rank 2, by Theorem 3 the parameters of these conditioned models are identifiable, up to the labeling of the states of the hidden variable. As these assumptions are generic conditions on the parameters of the original model, we can generically identify the parameters of the conditioned models.

In particular, M_4 can be identified up to reordering its rows, and is invertible. But let U denote the (unknown) 8×2 matrix with $U((i, j, k), \ell) = P(X_0 = \ell, X_1 = i, X_2 = j, X_3 = k)$. Then $P = UM_4$ has as its entries the observable distribution $P(X_1, X_2, X_3, X_4)$. Thus $U = PM_4^{-1}$ can be determined from P . Since U is the distribution of the induced model on X_0, X_1, X_2, X_3 with no hidden variables, it is then straightforward to identify all remaining parameters of the original model.

Thus all parameters are identifiable generically, up to label swapping. More specifically, they are identifiable provided that for either $j = 0$ or 1 the three matrices $M_4, M_2^{(j)}, M_3^{(j)}$ have rank 2, and \mathbf{p}_0 and the j th column of M_1 have non-zero entries.

5.3 Models 4-3e,f

Due to Markov equivalence, we need consider only 4-3e, as shown in Figure 6.

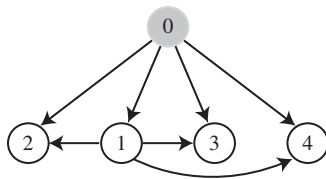


Figure 6 The DAG of model 4-3e

By conditioning on $X_1 = j, j = 1, 2$, we obtain two models of the form of 3-0. One checks that the induced parameters for these conditioned models are generic. Indeed, in terms of the original parameters they are $P(X_i|X_0, X_1 = j), i = 2, 3, 4$, which are generically non-singular since they are simply submatrices of $P(X_i|X_0, X_1)$, and at the hidden node

$$P(X_0|X_1 = j) = \frac{P(X_1 = j|X_0)P(X_0)}{\sum_{\ell} P(X_1 = j|X_0 = \ell)P(X_0 = \ell)},$$

which generically has non-zero entries.

Thus for generic parameters on the original model, up to label swapping we can determine $P(X_0|X_1 = j)$ and $P(X_i|X_0, X_1 = j), i = 2, 3, 4$. However, we do not have an ordering of the states of X_0 that is consistent for the recovered parameters for the two models. Generically we have four choices of parameters for the two models taken together. Each of these four choices leads to a possible joint distribution $P(X_0, X_1)$; viewing this joint distribution as a matrix, the four versions differ only by independently interchanging the two entries in each column, thus keeping the same marginalization $P(X_1)$. Generically, each of the four

distributions for $P(X_0, X_1)$ yields different parameters $P(X_0)$ and $P(X_1|X_0)$. The matrices $P(X_i|X_0, X_1)$, $i = 2, 3, 4$, are then obtained using the same rows as in $P(X_i|X_0, X_1 = j)$, though the ordering of the rows is dependent on the choice made previously.

Having obtained four possible parameter choices, it is straightforward to confirm that they all lead to the same joint distribution. Thus the parameterization map is generically 4-to-one.

6 Identification of causal effects

Here we examine the impact of k -to-one model parameterizations on the causal effect of one observable variable on another, when a latent variable acts as a confounder. For simplicity we assume that the latent variable is binary, though our discussion can be extended to a more general setting.

According to Theorem 3.2.2 (Adjusting for Direct Causes), p. 73 of Pearl [4], the causal effect of X_i on X_j can be obtained from model parameters by an appropriate sum over the states of the other direct causes of X_j . This sum is invariant under a relabeling of the states of those direct causes, and therefore the causal effect is not affected by label swapping if one of these is latent. As an instance, the causal effect of X_1 on X_2 in model 4-2b is

$$P(X_2|do(X_1 = x_1)) = P(X_2|X_1 = x_1, X_0 = 1)P(X_0 = 1) + P(X_2|X_1 = x_1, X_0 = 2)P(X_0 = 2). \tag{8}$$

Thus when label swapping is the only source of parameter non-identifiability, causal effects are uniquely determined by the observable distribution.

Things are more complex when parameter non-identifiability arises in other ways. For example, model 4-3e has one binary latent variable but a 4-to-one parameterization. In Table 1 two choices of parameters, (1), (2), are given for this model, as well as the common observable distribution they produce. These parameters and their two variants from label swapping at node 0 give the four elements of the fiber of the observable distribution.

Table 1 A rational example for model 4-3e. The parameter choices (1) and (2) lead to the same observable distribution, shown at the bottom. For the 4×2 matrix parameters, row indices refer to states of a pair of parents $i < j$ ordered lexicographically as (1,1), (1,2), (2,1), (2,2), with the first entry referring to parent i , and the second to parent j

(1)	$\mathbf{p}_0 = (2/5 \quad 3/5) \quad \mathbf{M}_1 = \begin{pmatrix} 2/5 & 3/5 \\ 14/15 & 1/15 \end{pmatrix}$
	$\mathbf{M}_2 = \begin{pmatrix} 2/5 & 3/5 \\ 3/5 & 2/5 \\ 4/5 & 1/5 \\ 9/10 & 1/10 \end{pmatrix} \quad \mathbf{M}_3 = \begin{pmatrix} 1/5 & 4/5 \\ 9/20 & 11/20 \\ 1/2 & 1/2 \\ 2/5 & 3/5 \end{pmatrix} \quad \mathbf{M}_4 = \begin{pmatrix} 1/2 & 1/2 \\ 7/10 & 3/10 \\ 4/5 & 1/5 \\ 3/5 & 2/5 \end{pmatrix}$
(2)	$\mathbf{p}'_0 = (1/5 \quad 4/5) \quad \mathbf{M}'_1 = \begin{pmatrix} 4/5 & 1/5 \\ 7/10 & 3/10 \end{pmatrix}$
	$\mathbf{M}'_2 = \begin{pmatrix} 2/5 & 3/5 \\ 9/10 & 1/10 \\ 4/5 & 1/5 \\ 3/5 & 2/5 \end{pmatrix} \quad \mathbf{M}'_3 = \begin{pmatrix} 1/5 & 4/5 \\ 2/5 & 3/5 \\ 1/2 & 1/2 \\ 9/20 & 11/20 \end{pmatrix} \quad \mathbf{M}'_4 = \begin{pmatrix} 1/2 & 1/2 \\ 3/5 & 2/5 \\ 4/5 & 1/5 \\ 7/10 & 3/10 \end{pmatrix}$
	$P(X_1, X_2, X_3 = 1, X_4 = 1) = \begin{bmatrix} 116/625 & 34/625 \\ 27/500 & 39/1250 \end{bmatrix}$
	$P(X_1, X_2, X_3 = 1, X_4 = 2) = \begin{bmatrix} 32/625 & 13/625 \\ 63/2500 & 17/1250 \end{bmatrix}$
	$P(X_1, X_2, X_3 = 2, X_4 = 1) = \begin{bmatrix} 128/625 & 52/625 \\ 171/2500 & 24/625 \end{bmatrix}$
	$P(X_1, X_2, X_3 = 2, X_4 = 2) = \begin{bmatrix} 44/625 & 31/625 \\ 81/2500 & 21/1250 \end{bmatrix}$

For any parameters of model 4-3e, the causal effect of X_1 on X_2 is again as given in eq. (8). However, due to the 4-to-one parameterization, there can be two different causal effects that are consistent with an observable distribution. As such, there may be distributions such that one causal effect leads to the conclusion that there is a positive effect of X_1 on X_2 (i.e., setting $X_1 = 2$ gives a higher probability of $X_2 = 2$ than setting $X_1 = 1$), while the other causal effect leads to the conclusion that there is a negative effect of X_1 on X_2 . Indeed, the observable distribution in Table 1 is such an instance. In Table 2 the two causal effects corresponding to that given distribution are shown. Here parameters (2) lead to a positive effect of X_1 on X_2 , while parameters (1) lead to a negative effect.

Table 2 The causal effects of X_1 on X_2 for the example in Table 1

Parameters	(a) $P(X_2 = 2 do(X_1 = 2))$	(b) $P(X_2 = 2 do(X_1 = 1))$	(a)–(b)
(1)	11/50	9/25	–7/50
(2)	17/50	7/25	3/50

More generally, for generic observable distributions of this model, choices of parameters that differ other than by label swapping can give different causal effects. However, it varies whether the effects have the same or different signs.

7 Conclusion

Paraphrasing Pearl [6], the problem of identifying causal effects in non-parametric models has been “placed to rest” by the proof of completeness of the *do*-calculus and related graphical criteria. In this paper we show that the introduction of modest (parametric) assumptions on the size of the state spaces of variables allows for identifiability of parameters that otherwise would be non-identifiable. Causal effects can be computed from identified parameters, if desired, but our techniques allow for the recovery of all parameters. In the process of proving parameter identifiability for several small networks, we use techniques inspired by a theorem of Kruskal, and other novel approaches. This framework can be applied to other models as well.

We have at least three reasons to extend the work described in this paper. The first is to develop new techniques and to prove new theoretical results for parameter identifiability; this provides the foundation of our work. The second is to reach the stage at which one can easily determine parameter identifiability for DAG models with hidden variables that are used in statistical modeling; this motivates our work. A third and related focus of future work is to address the scalability of our approach and to automate it. As noted above many of our arguments do not depend on variables being binary. Also, a strategy that we used successfully to handle larger models is to first marginalize or condition on an observable variable to reduce the model to one already understood, and then to “lift” results on the reduced model back to the original one. We are working toward turning this strategy into an algorithm.

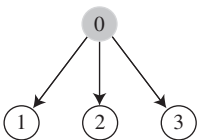
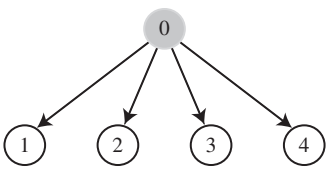
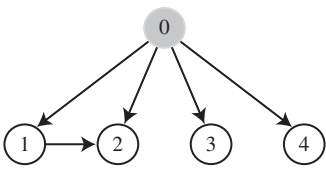
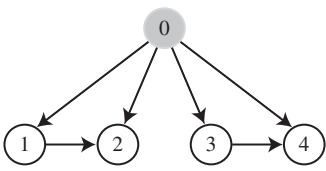
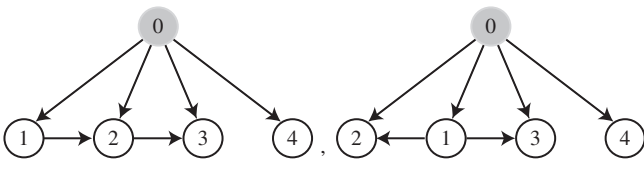
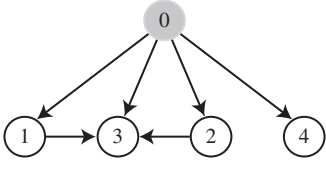
Acknowledgments: The authors thank the American Institute of Mathematics, where this work was begun during a workshop on Parameter Identification in Graphical Models, and continued through AIM’s SQuaRE program.

Research funding: American Institute of Mathematics Structured Quartet Research Ensemble grant.

Appendix

Table 3 shows all DAGs with four or fewer observable nodes and one hidden node that is a parent of all observable ones (see Section 5 for model naming convention). Markov equivalent graphs appear on the same line. The dimension of the parameter space is $\dim(\Theta)$, and $2^A - 1$ is the dimension of the probability simplex in which the joint distribution lies. The parameterization map is generically k -to-one.

Table 3 Small binary DAG models

Model	Graph	$\dim(\Theta)$	$2^A - 1$	k
$2-B, B \geq 0$		≥ 5	3	∞
3-0		7	7	2
$3-Bx, B \geq 1$		≥ 9	7	∞
4-0		9	15	2
4-1		11	15	2
4-2a		13	15	∞
4-2b,c		13	15	2
4-2d		15	15	2

(continued)

Table 3 (Continued)

Model	Graph	$\dim(\Theta)$	$2^A - 1$	k
4-3a,b		15	15	2
4-3c,d		17	15	∞
4-3e,f		15	15	4
4-3g		17	15	∞
4-3h		25	15	∞
4-3i		25	15	∞
4-Bx, $B \geq 4$		≥ 19	15	∞

References

1. Neapolitan RE. Probabilistic reasoning in expert systems: theory and algorithms. New York, NY: John Wiley and Sons, 1990.
2. Neapolitan RE. Learning Bayesian networks. Upper Saddle River, NJ: Pearson Prentice Hall, 2004.
3. Pearl J. Causal diagrams for empirical research. *Biometrika* 1995;82:669–710.
4. Pearl J. Causality: models, reasoning, and inference, 2nd ed. Cambridge: Cambridge University Press, 2009.
5. Huang Y, Valtorta M. Pearl's calculus of intervention is complete. In Proceedings of the twenty-second conference on uncertainty in artificial intelligence (UAI-06), 2006:217–24.
6. Pearl J. The do-calculus revisited. In Proceedings of the twenty-eighth conference on uncertainty in artificial intelligence (UAI-12), 2012:4–11.
7. Shpitser I, Pearl J. Complete identification methods for the causal hierarchy. *J Mach Learn Res* 2008;9:1941–79.
8. Tian J, Pearl J. A general identification condition for causal effects. In Proceedings of the eighteenth national conference on artificial intelligence (AAAI-02), 2002:567–73.
9. Kuroki M, Pearl J. Measurement bias and effect restoration in causal inference. *Biometrika* 2014;101:423–37.
10. Allman E, Matias C, Rhodes J. Identifiability of parameters in latent structure models with many observed variables. *Ann Statist* 2009;37:3099–132.
11. Mond D, Smith J, van Straten D. Stochastic factorizations, sandwiched simplices and the topology of the space of explanations. *R Soc Lond Proc Ser Math Phys Eng Sci* 2003;459:2821–45.
12. Kubjas K, Robeva E, Sturmfels B. Fixed points of the EM algorithm and nonnegative rank boundaries. In revision, 2013.
13. Kruskal J. Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra Appl* 1977;18:95–138.
14. Stanghellini E, Vantaggi B. Identification of discrete concentration graph models with one hidden binary variable. *Bernoulli* 2013;19:1920–37. Available at: <http://dx.doi.org/10.3150/12-BEJ435>.
15. Lauritzen SL. Graphical models. Oxford Statistical Science Series, vol. 17. New York: The Clarendon Press Oxford University Press, Oxford Science Publications, 1996.
16. Meek C. Strong completeness and faithfulness in Bayesian networks. In Proceedings of the eleventh annual conference on uncertainty in artificial intelligence (UAI-95), San Francisco, CA: Morgan Kaufmann, 1995:411–18.
17. Chickering DM. A transformational characterization of equivalent Bayesian network structures. In Proceedings of the eleventh annual conference on uncertainty in artificial intelligence (UAI-95), San Francisco, CA: Morgan Kaufmann, 1995:87–98.
18. Rhodes J. A concise proof of Kruskal's theorem on tensor decomposition. *Linear Algebra Appl* 2010;432:1818–24.
19. Stegeman A, Sidiropoulos ND. On Kruskal's uniqueness condition for the Candecomp/Parafac decomposition. *Linear Algebra Appl* 2007;420:540–52.