

Identifying the irreducible disjoint factors of a multivariate probability distribution.

Maxime Gasse, Alex Aussem

► **To cite this version:**

Maxime Gasse, Alex Aussem. Identifying the irreducible disjoint factors of a multivariate probability distribution.. Alessandro Antonucci; Giorgio Corani; Cassio Polpo de Campos. Probabilistic Graphical Models, Sep 2016, Lugano, Switzerland. JMLR.org, JMLR Workshop and Conference Proceedings, 52, pp.183 - 194, 2016, Proceedings of the Eighth International Conference on Probabilistic Graphical Models. <<http://jmlr.org/proceedings/papers/v52/>>. <hal-01425447>

HAL Id: hal-01425447

<https://hal.archives-ouvertes.fr/hal-01425447>

Submitted on 3 Jan 2017

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

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

Identifying the irreducible disjoint factors of a multivariate probability distribution

Maxime Gasse

Alex Aussem

University of Lyon, CNRS

Université Lyon 1, LIRIS, UMR5205, F-69622, France

MAXIME.GASSE@LIRIS.CNRS.FR

ALEXANDRE.AUSSEM@LIRIS.CNRS.FR

Abstract

We study the problem of decomposing a multivariate probability distribution $p(\mathbf{v})$ defined over a set of random variables $\mathbf{V} = \{V_1, \dots, V_n\}$ into a product of factors defined over disjoint subsets $\{\mathbf{V}_{F_1}, \dots, \mathbf{V}_{F_m}\}$. We show that the decomposition of \mathbf{V} into irreducible disjoint factors forms a unique partition, which corresponds to the connected components of a Bayesian or Markov network, given that it is faithful to p . Finally, we provide three generic procedures to identify these factors with $O(n^2)$ pairwise conditional independence tests $(V_i \perp\!\!\!\perp V_j | \mathbf{Z})$ under much less restrictive assumptions: i) p supports the Intersection property; ii) p supports the Composition property; iii) no assumption at all.

Keywords: conditional independence; probability distribution factorization; graphoids.

1. Introduction

The whole point of modeling a multivariate probability distribution p with a probabilistic graphical model, namely a Bayesian or a Markov network, is to encode independence relations into the graphical structure \mathcal{G} , thereby factorizing the joint probability distribution into a product of potential functions,

$$p(\mathbf{v}) = \prod_{i=1}^m \Phi_i(\mathbf{v}_i).$$

Such a factorization acts as a structural constraint on the expression of p , which reduces the number of free parameters in the model and facilitates both the learning and inference tasks, i.e. estimating p from a set of data samples, and answering probabilistic queries such as $\arg \max_{\mathbf{v}} p(\mathbf{v})$.

The fundamental problem that we wish to address in this paper involves finding a factorization of p into potential functions defined over minimal disjoint subsets, called *irreducible disjoint factors* (IDF as a shorthand). Such a factorization represents a strong structural constraint, and simplifies greatly the expression of p . For example, given two disjoint factors \mathbf{V}_1 and \mathbf{V}_2 , the task of obtaining $\arg \max_{\mathbf{v}} p(\mathbf{v})$ can be decomposed into two independent problems $\arg \max_{\mathbf{v}_1} p(\mathbf{v}_1)$ and $\arg \max_{\mathbf{v}_2} p(\mathbf{v}_2)$. Finding a set of disjoint factors is, for instance, an essential task in Sum-Product network (SPN) structure learning (Gens and Domingos, 2013), where product nodes correspond exactly to a product between disjoint factors, i.e. $p(\mathbf{v}_1) \times p(\mathbf{v}_2)$. Also, it was shown in (Gasse et al., 2015; Bielza et al., 2011) that identifying disjoint factors in a conditional distribution $p(\mathbf{y}|\mathbf{x})$ can effectively improve the maximum-a-posteriori estimation $\arg \max_{\mathbf{y}} p(\mathbf{y}|\mathbf{x})$ in multi-label classification.

In Section 2 we define our notations and introduce the concept of irreducible disjoint factors as well as some basic properties of conditional independence on which our theoretical results will

heavily rely. In Section 3 we show that irreducible disjoint factors necessarily form a unique partition, which relates to connected components in classical probabilistic graphical models. In Section 4, we establish several theoretical results to characterize the irreducible disjoint factors with pairwise conditional independence tests given several assumptions about p , namely the Intersection and Composition assumption, and then without any assumption. Each of these results establishes a quadratic generic procedure, which can be instantiated with only $O(n^2)$ statistical tests of independence. Finally, we conclude in Section 5.

2. Basic concepts

In this paper, upper-case letters in italics denote random variables (e.g. X, Y) and lower-case letters in italics denote their values (e.g. x, y). Likewise, upper-case bold letters denote random variable sets (e.g. $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$) and lower-case bold letters denote their values (e.g. $\mathbf{x}, \mathbf{y}, \mathbf{z}$). In the following we will consider only the multi-variate random variable $\mathbf{V} = \{V_1, \dots, V_n\}$ and its subsets. To keep the notation uncluttered, we use $p(\mathbf{v})$ to denote $p(\mathbf{V} = \mathbf{v})$ the joint distribution of \mathbf{V} . We recall the definition of conditional independence,

Definition 1 \mathbf{X} is conditionally independent of \mathbf{Y} given \mathbf{Z} , denoted $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}$, when $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subseteq \mathbf{V}$ are disjoint subsets of random variables such that for every value of $\mathbf{x}, \mathbf{y}, \mathbf{z}$ the following holds:

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z})p(\mathbf{z}) = p(\mathbf{x}, \mathbf{z})p(\mathbf{y}, \mathbf{z}).^1$$

We will assume the reader is familiar with the concept of separation in probabilistic graphical models, namely d -separation in directed acyclic graphs (DAGs) for Bayesian networks, and u -separation in undirected graphs (UGs) for Markov networks. These can be found in most books about probabilistic graphical models, e.g. Pearl (1989); Studeny (2005); Koller and Friedman (2009).

2.1 Disjoint factorization

We shall now introduce the concept of disjoint factors of random variables that will play a pivotal role in the factorization of the distribution $p(\mathbf{v})$.

Definition 2 A disjoint factor of random variables is a subset $\mathbf{V}_F \subseteq \mathbf{V}$ such that $\mathbf{V}_F \perp\!\!\!\perp \mathbf{V} \setminus \mathbf{V}_F$. Additionally, an irreducible disjoint factor is non-empty and has no other non-empty disjoint factor as proper subset.

As we will show next, irreducible disjoint factors necessarily form a partition of \mathbf{V} , which we denote \mathcal{F}_I . The key idea is then to decompose the joint distribution of the variables into a product of marginal distributions,

$$p(\mathbf{v}) = \prod_{\mathbf{V}_F \in \mathcal{F}_I} p(\mathbf{v}_F).$$

This paper aims to obtain theoretical results for the characterization of the irreducible disjoint factors with only pairwise conditional independence relations in the form $V_i \perp\!\!\!\perp V_j \mid \mathbf{Z}$.

1. Note that most definitions from the literature present the condition $p(\mathbf{x}, \mathbf{y} \mid \mathbf{z}) = p(\mathbf{x} \mid \mathbf{z})p(\mathbf{y} \mid \mathbf{z})$, which rely on the positivity condition $p(\mathbf{z}) > 0$.

2.2 Conditional independence properties

Consider four mutually disjoint random variables, \mathbf{W} , \mathbf{X} , \mathbf{Y} and \mathbf{Z} , and p the underlying probability distribution. As shown in (Dawid, 1979; Spohn, 1980), the properties of *Symmetry*, *Decomposition*, *Weak Union* and *Contraction* hold for any p , that is

$$\begin{aligned} \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} &\iff \mathbf{Y} \perp\!\!\!\perp \mathbf{X} \mid \mathbf{Z} \text{ (Symmetry),} \\ \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} &\implies \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \text{ (Decomposition),} \\ \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} &\implies \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \cup \mathbf{W} \text{ (Weak Union),} \\ \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \wedge \mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z} \cup \mathbf{Y} &\implies \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} \text{ (Contraction).} \end{aligned}$$

Any independence model that respects these four properties is called a *semi-graphoid* (Pearl and Verma, 1987). A fifth property holds in strictly positive distributions ($p > 0$), i.e. the *Intersection* property

$$\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \cup \mathbf{W} \wedge \mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z} \cup \mathbf{Y} \implies \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} \text{ (Intersection).}$$

Any independence model that respects these five properties is called a *graphoid*. The term "graphoid" was proposed by Pearl and Paz (1986) who noticed that these properties had striking similarities with vertex separation in graphs. Finally, a sixth property will be of particular interest in this work, that is the *Composition* property

$$\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \wedge \mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z} \implies \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} \text{ (Composition).}$$

The composition property holds in particular probability distributions, such as regular multivariate Gaussian distributions. Any independence model that respects these six properties is called a *compositional graphoid* (Sadeghi and Lauritzen, 2014). As shown in (Sadeghi and Lauritzen, 2015), independence models induced by classic probabilistic graphical models are compositional graphoids.

3. Problem analysis

Let us now develop further the notion of irreducible disjoint factors, and derive a first general graphical characterization. All proofs of the Theorems and Lemmas presented hereinafter are deferred to the Appendix.

3.1 Disjoint factors algebraic structure

We first show that disjoint factors can be characterized as an algebraic structure satisfying certain axioms. Let \mathcal{F} denote the set of all disjoint factors (DFs for short) defined over \mathbf{V} , and $\mathcal{F}_I \subset \mathcal{F}$ the set of all irreducible disjoint factors (IDFs for short). It is easily shown that $\{\mathbf{V}, \emptyset\} \subseteq \mathcal{F}$. More specifically, \mathcal{F} can be ordered via subset inclusion to obtain a lattice bounded by \mathbf{V} itself and the null set, while \mathcal{F}_I forms a partition of \mathbf{V} .

Theorem 3 *If $\mathbf{V}_{F_i}, \mathbf{V}_{F_j} \in \mathcal{F}$, then $\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j} \in \mathcal{F}$, $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j} \in \mathcal{F}$, and $\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j} \in \mathcal{F}$. Moreover, \mathbf{V} breaks down into a unique partition of irreducible components, \mathcal{F}_I .*

3.2 Graphical characterization

Irreducible disjoint factors will be conveniently represented as connected components in a graph, as we will see. Let us first introduce an important intermediary result.

Lemma 4 *Two distinct variables V_i and V_j belong to the same irreducible disjoint factor if there exists $\mathbf{Z} \subseteq \mathbf{V} \setminus \{V_i, V_j\}$ such that $V_i \perp\!\!\!\perp V_j \mid \mathbf{Z}$.*

Note that the converse is not true, so Lemma 4 does not provide a complete characterization of irreducible disjoint factors. The following example illustrates that point.

Example 1 *Consider $\mathbf{V} = \{V_1, V_2, V_3\}$, with V_3 a quaternary variable in $\{00, 01, 10, 11\}$ and V_1, V_2 two binary variables respectively equal to the first and the second digit of V_3 . Then, we have that $V_1 \perp\!\!\!\perp V_2$ and $V_1 \perp\!\!\!\perp V_2 \mid V_3$, and yet V_1 and V_2 belong to the same IDF $\{V_1, V_2, V_3\}$ due to $V_1 \not\perp\!\!\!\perp V_3$ and $V_2 \not\perp\!\!\!\perp V_3$.*

We now expand on Lemma 4 to propose a complete characterization of the irreducible disjoint factors using graph properties.

Theorem 5 *Let \mathcal{H} be an undirected graph whose nodes correspond to the random variables in \mathbf{V} , in which two nodes V_i and V_j are adjacent iff there exists $\mathbf{Z} \subseteq \mathbf{V} \setminus \{V_i, V_j\}$ such that $V_i \perp\!\!\!\perp V_j \mid \mathbf{Z}$. Then, each connected component in \mathcal{H} is an IDF.*

Theorem 5 offers an elegant graphical approach to characterize the IDFs, by mere inspection of the connected components in a graph. The problem of identifying all these connected components can be solved efficiently using a breadth-first search algorithm. Despite the desirable simplicity of this graphical characterization, deciding upon whether $\exists \mathbf{Z} \subseteq \mathbf{V} \setminus \{V_i, V_j\}$ such that $V_i \perp\!\!\!\perp V_j \mid \mathbf{Z}$ remains a challenging combinatorial problem, an exhaustive search for \mathbf{Z} being computationally infeasible even for moderate amounts of variables. Moreover, a second issue is that performing a statistical test of independence conditioned on a large \mathbf{Z} can become problematic; in the discrete case the sample size required for high-confidence grows exponentially in the size of the conditioning set. We show next that it is possible to overcome these limitations by considering restrictive assumptions regarding p .

3.3 IDFs and PGM structures

Note that, due to the d -separation criterion for DAGs and the u -separation criterion for UGs, it is possible to read off the IDFs directly from a Bayesian network or Markov network structure, given that it is faithful to p .

Corollary 6 *Let \mathcal{G} be a Bayesian or Markov network structure that is faithful to p . Then, two variables V_i and V_j belong to the same IDF iff there is a path between them in \mathcal{G} .*

Corollary 6 bridges the gap between the notion of irreducible disjoint factors and classical probabilistic graphical models. Still, the problem of structure learning for Bayesian and Markov networks is known to be NP-hard in general (Chickering et al., 2004; Karger and Srebro, 2001), and we have no guarantee that the probability distribution underlying the data is faithful to a Bayesian network or a Markov network structure. In the next section we consider practical procedures inspired from constraint-based structure learning algorithms, which allow us to extract the IDFs without relying on a particular PGM structure.

4. Generic procedures

In this section, we address the problem of identifying the irreducible disjoint factors from pairwise conditional independence tests. Finding a sound and efficient algorithmic procedure for general distributions is not completely trivial as we shall see, so we may consider several (reasonable) assumptions about the underlying distribution p , namely the Intersection and Composition properties.

4.1 Under the Intersection assumption

Let us present first a simplified characterization of the IDFs for distributions satisfying the Intersection property.

Theorem 7 *Let \mathcal{K} be an undirected graph whose nodes correspond to the random variables in \mathbf{V} , in which two nodes V_i and V_j are adjacent iff $V_i \not\perp\!\!\!\perp Y_j \mid \mathbf{V} \setminus \{V_i, V_j\}$. Then, each connected component in \mathcal{K} is an IDF if p supports the Intersection property.*

The graph, \mathcal{K} , is referred to as a *concentration graph* in the statistical literature (Cox and Wermuth, 1993). Theorem 7 is appealing compared to Theorem 5, as it greatly reduces computational expense incurred in obtaining the irreducible disjoint factors, with only a quadratic number of conditional independence tests. The graph \mathcal{K} may not be identical the graph \mathcal{H} from Theorem 5, though under the Intersection assumption their connected components are the same. Still, the size of the conditioning set $\mathbf{V} \setminus \{V_i, V_j\}$ is problematic for large variable sets, as it greatly reduces the confidence of a statistical test with limited samples. However, under the Intersection assumption the problem of performing that statistical test can be translated into a Markov boundary discovery problem, which can be solved with any off-the-shelf minimal feature subset selection algorithm.

Lemma 8 *Consider $V_i, V_j \in \mathbf{V}$ two distinct variables, and \mathbf{M}_i a Markov boundary of V_i in \mathbf{V} . Then, $V_j \notin \mathbf{M}_i$ implies $V_i \perp\!\!\!\perp V_j \mid \mathbf{V} \setminus \{V_i, V_j\}$, and the converse holds when p supports the Intersection property.*

Note that the Intersection assumption might be too restrictive in many practical scenarios. In fact, many real-life distributions (e.g. engineering systems such as digital circuits and engines that contain deterministic components) violate the Intersection property. As noted in (Statnikov et al., 2013), high-throughput molecular data, known as the “multiplicity” of molecular signatures (i.e., different gene/biomarker sets perform equally well in terms of predictive accuracy of phenotypes) also suggests existence of multiple Markov boundaries, which violates Intersection. It is usually unknown to what degree the Intersection assumption holds in distributions encountered in practice. The following example provides a particular case where the Intersection property does not hold.

Example 2 *Consider $\mathbf{V} = \{V_1, V_2, V_3, V_4\}$ four binary random variables such that $v_1 = v_2$, $v_3 = v_4$ and $p(V_2 = V_3) = \alpha$, $0.5 < \alpha < 1$. Clearly here p is not strictly positive and the Intersection property does not hold. If we apply Theorem 7 along with Lemma 8 then we have $\mathbf{M}_1 = \{V_2\}$, $\mathbf{M}_2 = \{V_1\}$, $\mathbf{M}_3 = \{V_4\}$ and $\mathbf{M}_4 = \{V_3\}$, which results in two connected components $\{V_1, V_2\}$ and $\{V_3, V_4\}$ in \mathcal{K} . These are clearly not disjoint factors since $\{V_1, V_2\} \not\perp\!\!\!\perp \{V_3, V_4\}$.*

4.2 Under the Composition assumption

Second, we consider an even simpler characterization of the IDFs for distributions satisfying the Composition property.

Theorem 9 *Let \mathcal{L} be an undirected graph whose nodes correspond to the random variables in \mathbf{V} , in which two nodes V_i and V_j are adjacent iff $V_i \not\perp V_j$. Then, each connected component in \mathcal{L} is an IDF if p supports the Composition property.*

The graph, \mathcal{L} , is referred to as a *covariance graph* in the statistical literature (Cox and Wermuth, 1993). Theorem 9 is very similar to Theorem 7, with again a quadratic number of conditional independence tests involved. The graph \mathcal{L} may not be identical the graph \mathcal{H} from Theorem 5, though under the Composition assumption their connected components are the same. Moreover, a desirable property of this characterization is that the conditioning set vanishes, which ensures high confidence when performing a statistical test from finite samples. Still, it is usually unknown to what degree the Composition assumption holds in distributions encountered in practice. Some special distributions are known to satisfy the Composition property, for example multivariate Gaussian distributions (Studeny, 2005, Corollary 2.4) and the symmetric binary distributions used in (Wermuth et al., 2009). The following example provides a case where the Composition property does not hold.

Example 3 *Consider $\mathbf{V} = \{V_1, V_2, V_3\}$ three binary variables such that V_2 and V_3 are independent and uniformly distributed, and $p(V_1 = V_2 \oplus V_3) = \alpha$, $0.5 < \alpha < 1$ (\oplus denotes the exclusive OR operator). If we apply Theorem 9 we have that every pair of variables is mutually independent, which results in three connected components $\{V_1\}$, $\{V_2\}$ and $\{V_3\}$ in \mathcal{L} . These are clearly not disjoint factors since $\{V_1\} \not\perp \{V_2, V_3\}$, $\{V_2\} \not\perp \{V_3, V_1\}$ and $\{V_3\} \not\perp \{V_1, V_2\}$.*

4.3 For any probability distribution

Finally, we present a less trivial characterization of the IDFs that also loosens the computational burden by orders of magnitude compared to Theorem 5, and yet does not require any assumption about p .

Theorem 10 *Consider \prec a strict total order of \mathbf{V} . Let \mathcal{M} be an undirected graph whose nodes correspond to the random variables in \mathbf{V} , obtained from the following procedure:*

- 1: $\mathcal{M} \leftarrow (\mathbf{V}, \emptyset)$ (empty graph)
- 2: **for all** $V_i \in \mathbf{V}$ **do**
- 3: $\mathbf{V}_{ind}^i \leftarrow \emptyset$
- 4: **for all** $V_j \in (V|V > V_i)$ (processed in \prec order) **do**
- 5: **if** $V_i \perp V_j \mid \{V|V < V_i\} \cup \mathbf{V}_{ind}^i$ **then**
- 6: $\mathbf{V}_{ind}^i \leftarrow \mathbf{V}_{ind}^i \cup \{V_j\}$
- 7: **else**
- 8: Insert a new edge (i, j) in \mathcal{M}

Then, each connected component in \mathcal{M} is an IDF.

Here again, the number of conditional independence tests required in Theorem 10 is quadratic in the number of variables. Compared to the previous results under the Intersection and the Composition properties, this new characterization has the desirable advantage of requiring no assumption

about the underlying distribution p . However, it suffers from two limitations: i) the conditioning set at line 5 ranges from \emptyset in the first iteration to $\mathbf{V} \setminus \{V_i, V_j\}$ in the last iteration, which is problematic in high-dimensional data; and ii) the whole procedure is prone to a propagation of error, since each iteration depends on the result of the previous tests to constitute the \mathbf{V}_{ind}^i set. Also, the procedure can not be fully run in parallel, contrary to the procedures in Theorems 7 and 9.

5. Discussion

We presented three procedures based on pairwise conditional independence tests to identify the irreducible disjoint factors of a multivariate probability distribution $p(\mathbf{v})$. These procedures require only a quadratic number of independence tests, between each pair of variables $V_i, V_j \in \mathbf{V}$. The first one is correct under the assumption that p supports the Intersection property, and involves conditional independence tests in the form $V_i \perp\!\!\!\perp V_j \mid \mathbf{V} \setminus \{V_i, V_j\}$. The second one is correct under the assumption that p supports the Composition property, and involves conditional independence tests in the form $V_i \perp\!\!\!\perp V_j \mid \emptyset$. Finally, the third procedure we propose is correct for any probability distribution p , and involves conditional independence tests in the form $V_i \perp\!\!\!\perp V_j \mid \mathbf{Z}$, where \mathbf{Z} is updated iteratively from the outcome of the previous tests, and ranges from \emptyset to $\mathbf{V} \setminus \{V_i, V_j\}$.

While these three procedures are mathematically sound, their effective implementation will necessarily rely on fallible statistical tests in order to decide upon $\{X\} \perp\!\!\!\perp \{Y\} \mid \mathbf{Z}$, that is, the rejection or acceptance of the null hypothesis of independence. Typically, these are based on either a G or a χ^2 test when the data set is discrete and a Fisher's Z test when it is continuous. These tests are likely to fail to decide on conditional dependence when the expected counts in the contingency table are small. In fact, the decision of accepting or rejecting the null hypothesis depends implicitly upon the degree of freedom of the test, which increases exponentially with the number of variables in the conditional set. The larger the size of the conditioning set, the less accurate are the conditional probability estimates and hence the less reliable are the independence tests. Practical solutions to deal with this problem typically involve permutation-based tests to estimate the *effective* degrees of freedom (Tsamardinos and Borboudakis, 2010), kernel-based tests (Zhang et al., 2011), or a combination of both (Doran et al., 2014).

Among the three generic procedures presented in Theorems 7, 9 and 10, the second procedure (under Composition) is the more appealing, in our view, since it relies on low-order conditional independence test, which are more robust in practice. Moreover, the Composition property is usually considered as a reasonable assumption, and often tacitly assumed. For example, linear models rely on the Composition property. In the context of feature subset selection, it is often argued that forward selection is computationally more efficient than backward elimination (Guyon and Elisseeff, 2003). In fact such a statement tacitly supposes that the Composition property holds (Peña et al., 2007). Interestingly, the procedure used for SPN structure learning in (Gens and Domingos, 2013) to "*partition \mathbf{V} into approximately independent subsets \mathbf{V}_j* " can be seen as a direct instantiation of Theorem 9 with a G test of pairwise independence. We proved therefore that this particular procedure is in fact correct and optimal (i.e., it yields independent and irreducible subsets) when p supports the Composition property. Finally, the problem of decomposing p into irreducible disjoint factors seems closely related to the so-called "all-relevant" feature subset selection problem discussed in (Rudnicki et al., 2015; Nilsson et al., 2007), where a variable V_i is said to be relevant to another variable V_j iff there exists $\mathbf{Z} \subseteq \mathbf{V} \setminus \{V_i, V_j\}$ such that $V_i \not\perp\!\!\!\perp V_j \mid \mathbf{Z}$. The graph in Theorem 5 provides a straightforward solution to this problem, therefore it may be interesting to investigate fur-

ther how the graphs in Theorems 7, 9 and 10 may solve the "all-relevant" feature selection problem. This is left for future work.

Acknowledgments

The authors are very grateful to the anonymous reviewers for their insightful comments. This work was funded by both the French state through the Nano 2017 investment program and the European Community through the European Nanoelectronics Initiative Advisory Council (ENIAC Joint Undertaking), under grant agreement no 324271 (ENI.237.1.B2013).

Appendix

For the sake of conciseness, the obvious Symmetry property (i.e., $\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}$ equivalent to $\mathbf{Y} \perp \mathbf{X} \mid \mathbf{Z}$) will be used implicitly in the proofs.

Proof [Proof of Theorem 3] First, we prove that $\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j} \in \mathcal{F}$. From the DF assumption for \mathbf{V}_{F_i} and \mathbf{V}_{F_j} we have $\mathbf{V}_{F_i} \perp \mathbf{V} \setminus \mathbf{V}_{F_i}$ and $\mathbf{V}_{F_j} \perp \mathbf{V} \setminus \mathbf{V}_{F_j}$. Using the Weak Union property we obtain that $\mathbf{V}_{F_i} \perp \mathbf{V} \setminus (\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j}) \mid \mathbf{V}_{F_j} \setminus \mathbf{V}_{F_i}$, and similarly with the Decomposition property we get $\mathbf{V}_{F_j} \setminus \mathbf{V}_{F_i} \perp \mathbf{V} \setminus (\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j})$. We may now apply the Contraction property to show that $\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j} \perp \mathbf{V} \setminus (\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j})$. Therefore, $\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j}$ is a DF by definition. Second, we prove that $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j} \in \mathcal{F}$. From the DF assumption for \mathbf{V}_{F_i} and \mathbf{V}_{F_j} we have $\mathbf{V}_{F_i} \perp \mathbf{V} \setminus \mathbf{V}_{F_i}$ and $\mathbf{V}_{F_j} \perp \mathbf{V} \setminus \mathbf{V}_{F_j}$. Using the Weak Union property we obtain $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j} \perp (\mathbf{V} \setminus (\mathbf{V}_{F_i} \cup \mathbf{V}_{F_j})) \cup (\mathbf{V}_{F_j} \setminus \mathbf{V}_{F_i}) \mid \mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j}$, and similarly with the Decomposition property we get $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j} \perp \mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j}$. We may now apply the Contraction property to show that $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j} \perp \mathbf{V} \setminus (\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j})$. Therefore, $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j}$ is a DF by definition. Third, we prove that $\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j} \in \mathcal{F}$. From the DF assumption for \mathbf{V}_{F_i} and \mathbf{V}_{F_j} we have $\mathbf{V}_{F_i} \perp \mathbf{V} \setminus \mathbf{V}_{F_i}$ and $\mathbf{V}_{F_j} \perp \mathbf{V} \setminus \mathbf{V}_{F_j}$. Using the Weak Union property we obtain $\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j} \perp \mathbf{V} \setminus \mathbf{V}_{F_i} \mid \mathbf{V}_{F_j}$, and similarly with the Decomposition property we get $\mathbf{V}_{F_j} \perp \mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j}$. We may now apply the Contraction property to show that $\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j} \perp \mathbf{V} \setminus (\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j})$. Therefore, $\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j}$ is a DF by definition. Finally, we prove that \mathcal{F}_I forms a partition of \mathbf{V} . Consider a non-empty DF $\mathbf{V}_{F_i} \in \mathcal{F}$. Then either \mathbf{V}_{F_i} is an IDF, or one of its proper non-empty subsets $\mathbf{V}_{F_j} \subset \mathbf{V}_{F_i}$ is an IDF and the remaining set $\mathbf{V}_{F_i} \setminus \mathbf{V}_{F_j}$ is a non-empty DF. By applying the same reasoning recursively, the non-empty DF $\mathbf{V} \in \mathcal{F}$ breaks down into an irreducible partition of IDFs. Now, consider two distinct IDFs $\mathbf{V}_{F_i}, \mathbf{V}_{F_j} \in \mathcal{F}_I$, then $\mathbf{V}_{F_i} \cap \mathbf{V}_{F_j}$ is a DF, which is necessarily empty due to the IDF assumption for \mathbf{V}_{F_i} or \mathbf{V}_{F_j} . As a result all IDFs are mutually disjoint, and \mathcal{F}_I forms a unique partition of \mathbf{V} . ■

Proof [Proof of Lemma 4] By contradiction, suppose V_i and V_j do not belong to the same IDF, and let \mathbf{V}_{F_i} denote the irreducible disjoint factor to which V_i belongs. From the DF definition we have $\mathbf{V}_{F_i} \perp \mathbf{V} \setminus \mathbf{V}_{F_i}$. Let \mathbf{Z} denote any arbitrary subset of $\mathbf{V} \setminus \{V_i, V_j\}$, we can apply the Weak Union property to obtain $\mathbf{V}_{F_i} \setminus \mathbf{Z} \perp \mathbf{V} \setminus (\mathbf{V}_{F_i} \cup \mathbf{Z}) \mid \mathbf{Z}$. Then, from the Decomposition property we have $\{V_i\} \perp \{V_j\} \mid \mathbf{Z}$. This is true for every such \mathbf{Z} subset, which concludes the proof. ■

We now introduce Lemma 11 which will prove useful to our subsequent demonstrations.

Lemma 11 *Let \mathbf{V}_F be an IDF. Then, for every nonempty proper subset \mathbf{Z} of \mathbf{V}_F , we have $\mathbf{Z} \not\perp\!\!\!\perp \mathbf{V}_F \setminus \mathbf{Z} \mid \mathbf{V} \setminus \mathbf{V}_F$.*

Proof [Proof of Lemma 11] By contradiction, suppose such a \mathbf{Z} exists with $\mathbf{Z} \perp\!\!\!\perp \mathbf{V}_F \setminus \mathbf{Z} \mid \mathbf{V} \setminus \mathbf{V}_F$. From the DF assumption of \mathbf{V}_F , we also have that $\mathbf{V}_F \perp\!\!\!\perp \mathbf{V} \setminus \mathbf{V}_F$, and therefore $\mathbf{Z} \perp\!\!\!\perp \mathbf{V} \setminus \mathbf{V}_F$ due to the Decomposition property. We may now apply the Contraction property on these two statements to obtain $\mathbf{Z} \perp\!\!\!\perp \mathbf{V} \setminus \mathbf{Z}$ which contradicts the IDF assumption for \mathbf{V}_F . This concludes the proof. ■

Proof [Proof of Theorem 5] If a path exists between V_i and V_j in \mathcal{H} then owing to Lemma 4 all pairs of successive variables in the path are in the same IDF, and by transitivity V_i and V_j necessarily belong to the same IDF. We may now prove the converse. Suppose that V_i and V_j belong to the same IDF, denoted \mathbf{V}_F . Consider $\{\mathbf{X}, \mathbf{Y}\}$ a partition of \mathbf{V} such that $V_i \in \mathbf{X}$ and $V_j \in \mathbf{Y}$. Then, owing to Lemma 11, we have that $\mathbf{X} \cap \mathbf{V}_F \not\perp\!\!\!\perp \mathbf{Y} \cap \mathbf{V}_F \mid \mathbf{V} \setminus \mathbf{V}_F$. Using the Weak Union property, we obtain $\mathbf{X} \not\perp\!\!\!\perp \mathbf{Y}$. Consider X_1 an arbitrary variable from \mathbf{X} . Using the Contraction property, we have that either $\{X_1\} \not\perp\!\!\!\perp \mathbf{Y}$ or $\mathbf{X} \setminus \{X_1\} \not\perp\!\!\!\perp \mathbf{Y} \mid \{X_1\}$. Consider X_2 another arbitrary variable from $\mathbf{X} \setminus \{X_1\}$, we can apply the Contraction property again on the second expression to obtain that either $\{X_2\} \not\perp\!\!\!\perp \mathbf{Y} \mid \{X_1\}$ or $\mathbf{X} \setminus \{X_1, X_2\} \not\perp\!\!\!\perp \mathbf{Y} \mid \{X_1, X_2\}$. If we proceed recursively, we will necessarily find a variable $X_k \in \mathbf{X}$ such that $\{X_k\} \not\perp\!\!\!\perp \mathbf{Y} \mid \{X_1, \dots, X_{k-1}\}$. Likewise, we can proceed along the same line to exhibit a variable $Y_l \in \mathbf{Y}$ such that $\{X_k\} \not\perp\!\!\!\perp \{Y_l\} \mid \{X_1, \dots, X_{k-1}\} \cup \{Y_1, \dots, Y_{l-1}\}$. In other words, for every partition $\{\mathbf{X}, \mathbf{Y}\}$ of \mathbf{V} such that $V_i \in \mathbf{X}$ and $V_j \in \mathbf{Y}$, there exists at least one variable X in \mathbf{X} , one variable Y in \mathbf{Y} and one subset $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X, Y\}$, such that $\{X\} \not\perp\!\!\!\perp \{Y\} \mid \mathbf{Z}$. So there necessarily exists a path between V_i and V_j in \mathcal{H} . This concludes the proof. ■

Proof [Proof of Corollary 6] From the d -separation criterion in DAGs (resp. the u -separation criterion in UGs), if \mathcal{G} is faithful to p , then $\{V_i\} \not\perp\!\!\!\perp \{V_j\} \mid \mathbf{Z}$ iff there is an open path between V_i and V_j given \mathbf{Z} . Since every path can be made open by conditioning on its collider nodes (resp. on the empty set), then for every pair of distinct variables $V_i, V_j \in \mathbf{V}$ connected by a path, there exists a subset $\mathbf{Z} \subseteq \mathbf{V} \setminus \{V_i, V_j\}$ such that $V_i \not\perp\!\!\!\perp V_j \mid \mathbf{Z}$. Conversely, if there exists no path between V_i and V_j , then $\{V_i\} \perp\!\!\!\perp \{V_j\} \mid \mathbf{Z}$ for every such a \mathbf{Z} subset. This concludes the proof. ■

Proof [Proof of Theorem 7] If a path exists between V_i and V_j in \mathcal{K} then owing to Lemma 4 all pairs of successive variables in the path are in the same IDF, and by transitivity V_i and V_j necessarily belong to the same IDF. We may now prove the converse. Suppose that V_i and V_j belong to the same IDF, denoted \mathbf{V}_F . Consider $\{\mathbf{X}, \mathbf{Y}\}$ a partition of \mathbf{V} such that $V_i \in \mathbf{X}$ and $V_j \in \mathbf{Y}$. Then, owing to Lemma 11, we have that $\mathbf{X} \cap \mathbf{V}_F \not\perp\!\!\!\perp \mathbf{Y} \cap \mathbf{V}_F \mid \mathbf{V} \setminus \mathbf{V}_F$. Using the Weak Union property, we obtain $\mathbf{X} \not\perp\!\!\!\perp \mathbf{Y}$. Consider X_1 an arbitrary variable from \mathbf{X} . Using the Intersection property, we have that either $\{X_1\} \not\perp\!\!\!\perp \mathbf{Y} \mid \mathbf{X} \setminus \{X_1\}$ or $\mathbf{X} \setminus \{X_1\} \not\perp\!\!\!\perp \mathbf{Y} \mid \{X_1\}$. Consider X_2 another arbitrary variable from $\mathbf{X} \setminus \{X_1\}$, we can apply the Intersection property again on the second expression to obtain that either $\{X_2\} \not\perp\!\!\!\perp \mathbf{Y} \mid \mathbf{X} \setminus \{X_2\}$ or $\mathbf{X} \setminus \{X_1, X_2\} \not\perp\!\!\!\perp \mathbf{Y} \mid \{X_1, X_2\}$. If we proceed recursively, we will necessarily find a variable $X_k \in \mathbf{X}$ such that $\{X_k\} \not\perp\!\!\!\perp \mathbf{Y} \mid \mathbf{X} \setminus \{X_k\}$. Likewise, we can proceed along the same line to exhibit a variable $Y_l \in \mathbf{Y}$ such that $\{X_k\} \not\perp\!\!\!\perp \{Y_l\} \mid \mathbf{V} \setminus \{X_k, Y_l\}$. In other words, for every partition $\{\mathbf{X}, \mathbf{Y}\}$ of \mathbf{V} such that $V_i \in \mathbf{X}$ and $V_j \in \mathbf{Y}$, there exists at least one variable X in \mathbf{X} and one variable Y in \mathbf{Y} such that $\{X\} \not\perp\!\!\!\perp \{Y\} \mid \mathbf{V} \setminus \{X, Y\}$. So there necessarily

exists a path between V_i and V_j in \mathcal{K} . This concludes the proof. \blacksquare

Proof [Proof of Lemma 8] First, if $\{V_i\} \not\perp\!\!\!\perp \{V_j\} \mid \mathbf{V} \setminus \{V_i, V_j\}$ then from the Weak Union property we have that $\{V_i\} \not\perp\!\!\!\perp \mathbf{V} \setminus (\{V_i\} \cup \mathbf{Z}) \mid \mathbf{Z}$ for every $\mathbf{Z} \subseteq \mathbf{V} \setminus \{V_i, V_j\}$. In other words, there exists no Markov blanket (neither Markov boundary) of V_i in \mathbf{V} which does not contain V_j . Second, we prove the converse. Suppose the Markov boundary \mathbf{M}_i of V_i in \mathbf{V} contains V_j , then we have $\{V_i\} \not\perp\!\!\!\perp \{V_j\} \cup \mathbf{V} \setminus (\{V_i\} \cup \mathbf{M}_i) \mid \mathbf{M}_i \setminus \{V_j\}$. We can apply the Intersection property to obtain that either $\{V_i\} \not\perp\!\!\!\perp \{V_j\} \mid \mathbf{V} \setminus \{V_i, V_j\}$ or $\{V_i\} \not\perp\!\!\!\perp \mathbf{V} \setminus (\{V_i\} \cup \mathbf{M}_i) \mid \mathbf{M}_i$. The second statement contradicts the Markov blanket assumption for \mathbf{M}_i , so we necessarily have that $\{V_i\} \not\perp\!\!\!\perp \{V_j\} \mid \mathbf{V} \setminus \{V_i, V_j\}$. This concludes the proof. \blacksquare

Proof [Proof of Theorem 9] If a path exists between V_i and V_j in \mathcal{L} then owing to Lemma 4 all pairs of successive variables in the path are in the same IDF, and by transitivity V_i and V_j necessarily belong to the same IDF. We may now prove the converse. Suppose that V_i and V_j belong to the same IDF, denoted \mathbf{V}_F . Consider $\{\mathbf{X}, \mathbf{Y}\}$ a partition of \mathbf{V} such that $V_i \in \mathbf{X}$ and $V_j \in \mathbf{Y}$. Then, owing to Lemma 11, we have that $\mathbf{X} \cap \mathbf{V}_F \not\perp\!\!\!\perp \mathbf{Y} \cap \mathbf{V}_F \mid \mathbf{V} \setminus \mathbf{V}_F$. Using the Weak Union property, we obtain $\mathbf{X} \not\perp\!\!\!\perp \mathbf{Y}$. Consider X_1 an arbitrary variable from \mathbf{X} . Using the Composition property, we have that either $\{X_1\} \not\perp\!\!\!\perp \mathbf{Y}$ or $\mathbf{X} \setminus \{X_1\} \not\perp\!\!\!\perp \mathbf{Y}$. Consider X_2 another arbitrary variable from $\mathbf{X} \setminus \{X_1\}$, we can apply the Composition property again on the second expression to obtain that either $\{X_2\} \not\perp\!\!\!\perp \mathbf{Y}$ or $\mathbf{X} \setminus \{X_1, X_2\} \not\perp\!\!\!\perp \mathbf{Y}$. If we proceed recursively, we will necessarily find a variable $X_k \in \mathbf{X}$ such that $\{X_k\} \not\perp\!\!\!\perp \mathbf{Y}$. Likewise, we can proceed along the same line to exhibit a variable $Y_l \in \mathbf{Y}$ such that $\{X_k\} \not\perp\!\!\!\perp \{Y_l\}$. In other words, for every partition $\{\mathbf{X}, \mathbf{Y}\}$ of \mathbf{V} such that $V_i \in \mathbf{X}$ and $V_j \in \mathbf{Y}$, there exists at least one variable X in \mathbf{X} and one variable Y in \mathbf{Y} such that $\{X\} \not\perp\!\!\!\perp \{Y\}$. So there necessarily exists a path between V_i and V_j in \mathcal{L} . This concludes the proof. \blacksquare

Proof [Proof of Theorem 10] To keep the subsequent developments uncluttered, we consider without loss of generality that the variable set $\mathbf{V} = \{V_1, \dots, V_n\}$ is ordered according to $<$, so that $V_i < V_j \iff i < j$. Second, we denote $\mathbf{V}_{ind}^{i,j}$ the set \mathbf{V}_{ind}^i in its intermediary state at line 5 when V_j is being processed, while \mathbf{V}_{ind}^i denotes its state at the end of the procedure. Last, we adopt the notation $\{X \mid X > V_k\}$ and $\{Y \mid Y > V_k\}$ to denote respectively the sets $\{V \mid V > V_k, V \in \mathbf{X}\}$ and $\{V \mid V > V_k, V \in \mathbf{Y}\}$ (with \mathbf{X}, \mathbf{Y} subsets of \mathbf{V}), so that $X \in \mathbf{X}$ and $Y \in \mathbf{Y}$ by convention.

We start by proving that V_i and V_j are in the same IDF if V_i and V_j are connected in \mathcal{M} . If two variables V_p and V_q (with $V_p < V_q$) are adjacent in \mathcal{M} , then there exists a set $\mathbf{V}_{ind}^{p,q}$ such that $\{V_p\} \not\perp\!\!\!\perp \{V_q\} \mid \{V \mid V < V_p\} \cup \mathbf{V}_{ind}^{p,q}$, and from Lemma 4 V_p and V_q belong to the same IDF. Now, if a path exists between V_i and V_j in \mathcal{M} , then all pairs of successive variables in the path are in the same IDF, and by transitivity V_i and V_j belong to the same IDF.

To show the converse, we shall prove by contradiction that if V_i and V_j belong to the same IDF, then there exists a path between V_i and V_j in \mathcal{M} . Suppose there is no such path, then there exists a partition $\{\mathbf{X}, \mathbf{Y}\}$ of \mathbf{V} such that $V_i \in \mathbf{X}$, $V_j \in \mathbf{Y}$, and every variable in \mathbf{X} is non-adjacent to every variable in \mathbf{Y} . Equivalently, for every variable $V_k \in \mathbf{V}$ we have $\{X \mid X > V_k\} \subseteq \mathbf{V}_{ind}^k$ if $V_k \in \mathbf{Y}$, and $\{Y \mid Y > V_k\} \subseteq \mathbf{V}_{ind}^k$ if $V_k \in \mathbf{X}$. To proceed, we shall first prove by induction that

$$\forall k > i, \quad \{V_i\} \perp\!\!\!\perp \mathbf{V}_{ind}^{i,k} \mid \{V \mid V < V_i\}.$$

For $k = i + 1$, we have that $\mathbf{V}_{ind}^{i,k} = \emptyset$ so the result holds trivially. Suppose that $\{V_i\} \perp\!\!\!\perp \mathbf{V}_{ind}^{i,k} \mid \{V \mid V < V_i\}$ holds for some k . If $\{V_i\} \perp\!\!\!\perp \{V_k\} \mid \{V \mid V < V_i\} \cup \mathbf{V}_{ind}^{i,k}$, then $\mathbf{V}_{ind}^{i,k+1} = \mathbf{V}_{ind}^{i,k} \cup \{V_k\}$ and $\{V_i\} \perp\!\!\!\perp \mathbf{V}_{ind}^{i,k+1} \mid \{V \mid V < V_i\}$ due to the Contraction property. Otherwise, $\mathbf{V}_{ind}^{i,k+1} = \mathbf{V}_{ind}^{i,k}$ and we end up with the same result. Therefore, the result holds for every $k > i$ by induction, and setting $k = n$ yields $\{V_i\} \perp\!\!\!\perp \mathbf{V}_{ind}^i \mid \{V \mid V < V_i\}$. Now, we prove a second result by induction:

$$\forall k, \{X \mid X \geq V_k\} \perp\!\!\!\perp \{Y \mid Y \geq V_k\} \mid \{V \mid V < V_k\}.$$

For $k = n$, we have $\{X \mid X \geq V_k\} = \{Y \mid Y \geq V_k\} = \emptyset$ so the result holds trivially. Consider the previous variable, V_{k-1} , and suppose it belongs to \mathbf{X} , then due to our previous result we have $\{V_{k-1}\} \perp\!\!\!\perp \mathbf{V}_{ind}^{k-1} \mid \{V \mid V < V_{k-1}\}$. Since $\{Y \mid Y > V_{k-1}\} \subseteq \mathbf{V}_{ind}^{k-1}$, we may apply the Decomposition property to obtain $\{V_{k-1}\} \perp\!\!\!\perp \{Y \mid Y \geq V_{k-1}\} \mid \{V \mid V < V_{k-1}\}$. Combining that last expression with $\{X \mid X \geq V_k\} \perp\!\!\!\perp \{Y \mid Y \geq V_k\} \mid \{V \mid V < V_k\}$ yields $\{X \mid X \geq V_{k-1}\} \perp\!\!\!\perp \{Y \mid Y \geq V_{k-1}\} \mid \{V \mid V < V_{k-1}\}$ due to the Contraction property. The same demonstration holds if $V_{k-1} \in \mathbf{Y}$. Therefore, the result holds for every k by induction. Setting $k = 1$ in the expression above yields $\mathbf{X} \perp\!\!\!\perp \mathbf{Y}$, therefore V_i and V_j belong to distinct IDFs. This concludes the proof. \blacksquare

References

- C. Bielza, G. Li, and P. Larrañaga. Multi-dimensional classification with Bayesian networks. *International Journal of Approximate Reasoning*, 52(6):705–727, 2011.
- D. M. Chickering, D. Heckerman, and C. Meek. Large-Sample Learning of Bayesian Networks is NP-Hard. *Journal of Machine Learning Research*, 5:1287–1330, 2004.
- D. R. Cox and N. Wermuth. Linear dependencies represented by chain graphs. *Statistical Science*, 8(3):204–218, 1993. doi: 10.1214/ss/1177010887.
- A. P. Dawid. Conditional Independence in Statistical Theory. *Journal of the Royal Statistical Society, Series B*, 41:1–31, 1979.
- G. Doran, K. Muandet, K. Zhang, and B. Schölkopf. A permutation-based kernel conditional independence test. In N. L. Zhang and J. Tian, editors, *UAI*, pages 132–141. AUAI Press, 2014. ISBN 978-0-9749039-1-0.
- M. Gasse, A. Aussem, and H. Elghazel. On the Optimality of Multi-Label Classification under Subset Zero-One Loss for Distributions Satisfying the Composition Property. In *ICML*, volume 37 of *JMLR Proceedings*, pages 2531–2539, 2015.
- R. Gens and P. M. Domingos. Learning the Structure of Sum-Product Networks. In *ICML*, volume 28 of *JMLR Proceedings*, pages 873–880, 2013.
- I. Guyon and A. Elisseeff. An Introduction to Variable and Feature Selection. *Journal of Machine Learning Research*, 3:1157–1182, 2003.
- D. R. Karger and N. Srebro. Learning markov networks: maximum bounded tree-width graphs. In S. R. Kosaraju, editor, *SODA*, pages 392–401. ACM/SIAM, 2001. ISBN 0-89871-490-7.

- D. Koller and N. Friedman. *Probabilistic Graphical Models - Principles and Techniques*. MIT Press, 2009. ISBN 978-0-262-01319-2.
- R. Nilsson, J. M. Peña, J. Björkegren, and J. Tegnér. Consistent Feature Selection for Pattern Recognition in Polynomial Time. *Journal of Machine Learning Research*, 8:589–612, 2007.
- J. Pearl. *Probabilistic reasoning in intelligent systems - networks of plausible inference*. Morgan Kaufmann, 1989.
- J. Pearl and A. Paz. Graphoids: Graph-Based Logic for Reasoning about Relevance Relations or When would x tell you more about y if you already know z? In *ECAI*, pages 357–363, 1986.
- J. Pearl and T. Verma. The Logic of Representing Dependencies by Directed Graphs. In *AAAI*, pages 374–379, 1987.
- J. M. Peña, R. Nilsson, J. Björkegren, and J. Tegnér. Towards scalable and data efficient learning of Markov boundaries. *International Journal of Approximate Reasoning*, 45(2):211–232, 2007.
- W. R. Rudnicki, M. Wrzesien, and W. Paja. All Relevant Feature Selection Methods and Applications. In U. Stanczyk and L. C. Jain, editors, *Feature Selection for Data and Pattern Recognition*, volume 584 of *Studies in Computational Intelligence*, pages 11–28. Springer, 2015. ISBN 978-3-662-45619-4.
- K. Sadeghi and S. Lauritzen. Markov properties for mixed graphs. *Bernoulli*, 20(2):676–696, 2014. doi: 10.3150/12-BEJ502.
- K. Sadeghi and S. Lauritzen. Unifying Markov properties in graphical models. Unpublished manuscript, 2015.
- W. Spohn. Stochastic independence, causal independence, and shieldability. *Journal of Philosophical Logic*, 9(1):73–99, 1980.
- A. R. Statnikov, J. Lemeire, and C. F. Aliferis. Algorithms for discovery of multiple Markov boundaries. *Journal of Machine Learning Research*, 14(1):499–566, 2013.
- M. Studeny. *Probabilistic Conditional Independence Structures*. Springer, 2005. ISBN 978-1-84628-083-2.
- I. Tsamardinos and G. Borboudakis. Permutation testing improves bayesian network learning. In J. L. Balcázar, F. Bonchi, A. Gionis, and M. Sebag, editors, *ECML/PKDD (3)*, volume 6323 of *Lecture Notes in Computer Science*, pages 322–337. Springer, 2010. ISBN 978-3-642-15938-1.
- N. Wermuth, G. M. Marchetti, and D. R. Cox. Triangular systems for symmetric binary variables. *Electronic Journal of Statistics*, 3:932–955, 2009. doi: 10.1214/09-EJS439.
- K. Zhang, J. Peters, D. Janzing, and B. Schölkopf. Kernel-based conditional independence test and application in causal discovery. In F. G. Cozman and A. Pfeffer, editors, *UAI*, pages 804–813. AUAI Press, 2011. ISBN 978-0-9749039-7-2.