

# Causality in the Sciences

**Phyllis McKay Illari**

University of Kent

**Federica Russo**

Université catholique de Louvain, University of Kent

**Jon Williamson**

University of Kent

CLARENDON PRESS • OXFORD

2009



## WHEN ARE GRAPHICAL CAUSAL MODELS NOT GOOD MODELS?

**Jan Lemeire, Kris Steenhaut, Abdellah Touhafi**  
**Dept. of Electronics and Informatics (ETRO), Vrije Universiteit Brussel (VUB) -**  
**Interdisciplinary Institute for Broadband Technology (IBBT)**

### Abstract

The principle of Kolmogorov Minimal Sufficient Statistic (KMSS) states that the meaningful information of data is given by the regularities in the data. The KMSS is the minimal model that describes the regularities. The meaningful information given by a Bayesian network is the Directed Acyclic Graph (DAG) which describes a decomposition of the joint probability distribution into Conditional Probability Distributions (CPDs). If the description given by the Bayesian network is incompressible, the DAG is the KMSS and is faithful. We prove that if a faithful Bayesian network exists, it is the minimal Bayesian network. Moreover, if a Bayesian network gives the KMSS, modularity of the CPDs is the most plausible hypothesis, from which the causal interpretation follows. On the other hand, if the minimal Bayesian network is compressible and is thus not the KMSS, the above implications cannot be guaranteed. When the non-minimality of the description is due to the compressibility of an individual CPD, the true causal model is an element of the set of minimal Bayesian networks and modularity is still plausible. Faithfulness cannot be guaranteed though. When the concatenation of the descriptions of the CPDs is compressible, the true causal model is not necessarily an element of the set of minimal Bayesian networks. Also modularity may become implausible. This suggests that either there is a kind of meta-mechanism governing some of the mechanisms or a wrong model class is considered.

### 1.1 Introduction

Inductive inference comes to modeling the patterns in the data. Patterns or regularities in observations are - most likely - not coincidences, but give us valuable information about the system under study. A regularity is identified by its ability to compress the data, i.e. to describe the data using fewer symbols than the number of symbols needed to describe the data literally. Compressiveness is objectively defined by the Kolmogorov complexity. The concept is, however, not directly applicable since there does not exist an algorithm that computes the shortest program for a string. Kolmogorov complexity is therefore mainly used for giving preference within a given set of models. This has given rise to different methods for inductive inference, such as Minimum Message Length and Minimum Description Length. These methods are used for selecting the best model

from a given set of models, the model class. The choice of model class, however, determines the regularities under consideration.

For analyzing the validity of causal inference, we do not want to stick to an a priori chosen set of regularities, but search for *all* relevant regularities. This idea is captured by the concept of Kolmogorov Minimal Sufficient Statistic (KMSS). The KMSS is the minimal model such that the model together with the data is described minimally. The model should capture all regularities and nothing more.

For causal inference, the set of Bayesian networks is used as a model class. The DAG of a Bayesian network gives a minimal description of the conditional independencies following from a causal structure. A system can, however, contain other regularities. Then, the assumptions and implications of causal model theory, such as faithfulness, modularity and the correctness of causal inference, may become invalid. It can give rise to other independencies so that the DAG becomes unfaithful. We will show that the presence of other regularities cannot be ignored.

In Section 2, we will introduce the concept of KMSS. In Section 3, we will give a survey of causal model theory and the learning algorithms. Section 4 discusses related work. In Section 5 we apply the principle of KMSS to inductive inference and show that a Bayesian network captures dependencies between variables. Section 6 establishes the link between minimality of Bayesian networks, compressibility and faithfulness. In Section 7 we will argue that causal inference is plausible if the minimal Bayesian network is the KMSS. Section 8 discusses various cases in which the minimal Bayesian network does not provide the minimal description.

## 1.2 Meaningful Information

The *Kolmogorov Complexity* of a string  $x$  is defined to be the length of the shortest computer program that prints the string and then halts (Li and Vitányi, 1997):

$$K(x) = \min_{p: \mathcal{U}(p)=x} l(p) \quad (1.1)$$

with  $\mathcal{U}$  a universal computer and  $l(p)$  the size in bits of program  $p$ . Patterns in the string allow for its compression, i.e. to describe the data using fewer symbols than the number of symbols to describe the data literally.

The string "000100010001000100010001000100010001000100010001" can be described shorter by program REPEAT 11 TIMES "0001". But not all bits of this program can be regarded as containing *meaningful information*. We consider meaningful information as the properties of the string that allow for its compression (Vitányi, 2002). Such properties are called patterns or *regularities*. The regularity of the example string is the repetition. The number of repetitions ("11") or the substring "0001" is random information. A random string, which is incompressible has no meaningful information at all.

For inductive inference, we will look for a minimal description in 2 parts, one containing the regularities or patterns of the data, which we put in the model, and one part containing the remaining random noise. Such a description is called a *two-part code*. This results in a generic approach for inductive inference, called *Minimum Description*

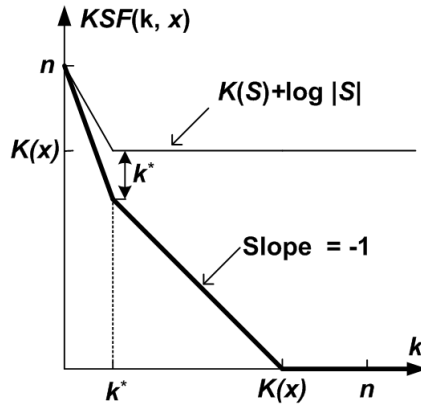


FIG. 1.1. Kolmogorov structure function for  $n$ -bit string  $x$ ,  $k^*$  is the KMSS of  $x$ .

*Length* (MDL), according to which we have to pick out the model  $M_{mdl}$  from model class  $\mathcal{M}$  where  $M_{mdl}$  is the model which minimizes the sum of the description length of  $M$  and of the data  $D$  encoded with the help of  $M$  (Grünwald, 1998):

$$M_{mdl} = \arg \min_{M \in \mathcal{M}} \{L(M) + L(D | M)\} \quad (1.2)$$

with  $L(\cdot)$  the description length.

The MDL approach relies on the a priori chosen model class. It does not tell us how to make sure the models capture all regularities of the data. The KMSS provides a formal separation of meaningful and meaningless information. We limit the introduction of KMSS to models that can be related to a finite set of objects, called the *model set*. In the context of learning, we are interested in a model set  $S$  that contains string  $x$  and the objects that share  $x$ 's regularities. All elements of a set  $S$  can be enumerated with a binary index of length  $\log_2 |S|$  with  $|S|$  the size of set  $S$ . We therefore say that  $x$  is *typical* for  $S$  if

$$K(x | p_S) \geq \log_2 |S| - \beta \quad (1.3)$$

with  $K(x | p_S)$  the conditional Kolmogorov complexity.  $p_S$  denotes the shortest program that describes  $S$  and  $\beta$  an agreed upon constant. Given set  $S$ ,  $x$  cannot be described shorter than by the set's index. Atypical elements have regularities that are not shared by most of the set's members and can therefore be described by a shorter description. Note that most elements of a set are typical, since, by counting arguments, only a small portion of it can be described shorter than  $\log_2 |S|$ .

The construction of  $S$  can be understood with the *Kolmogorov structure function*  $KSF$ .  $KSF(k, x)$  of  $x$  is defined as the  $\log_2$ -size of the smallest set including  $x$  which can be described with no more than  $k$  bits (Cover and Thomas, 1991):

$$KSF(k, x) = \min_{\substack{p: l(p) \leq k \\ \mathcal{U}(p) = S \\ x \in S}} \log_2 |S| \quad (1.4)$$

A typical graph of the structure function is illustrated in Figure 1.1. By taking  $k = 0$ , the only set that can be described is the entire set  $\{0, 1\}^n$  containing  $2^n$  elements, so that the corresponding log set size is  $n$ . By increasing  $k$ , the model can take advantage of the regularities of  $x$  in such way that each bit reduces the set's size more than halving it. The slope of the curve is smaller than -1. When  $k$  reaches  $k^*$ , all regularities are exploited. There are no more patterns in the data that allow for further compression. From then on each additional bit of  $k$  reduces the set by half. We proceed along the line of slope -1 until  $k = K(x)$  and the smallest set that can be described is the singleton  $\{x\}$ . The curve  $K(S) + \log_2 |S|$  is also shown on the graph. It represents the descriptive complexity of  $x$  by using the two-part code. With  $k = k^*$  it reaches its minimum and equals to  $K(x)$ . When  $k < k^*$ ,  $S$  is too general and is not a typical set for  $x$ .  $x$  is only typical for  $S$  if  $k \geq k^*$ . For random strings the curve starts at  $\log_2 |S| = n$  for  $k=0$  and drops with a slope of -1 until reaching the x-axis at  $k = n$ . Each bit reveals one of the bits of  $x$ , and halves the model set.

The *Kolmogorov Minimal Sufficient Statistic* (KMSS) of  $x$  is defined as the shortest program  $p^*$  which describes the smallest set  $S^*$  such that the two-stage description of  $x$  is as good as the minimal single-stage description of  $x$  (Gács et al., 2001; Vitányi, 2002):

$$p^* = \arg \min_p \{l(p) \mid \mathcal{U}(p) = S^*, x \in S^*, K(S^*) + \log_2 |S^*| \leq K(x)\} \quad (1.5)$$

The descriptive complexity of  $S^*$  is then  $k^*$ . Program  $p^*$  minimally describes the meaningful information present in  $x$  and nothing else. The definition ensures that  $x$  is a typical element of  $S^*$ .

### 1.3 Graphical Causal Models

This chapter will introduce graphical causal models and the accompanying learning algorithms (Pearl, 2000; Spirtes et al., 1993).

#### 1.3.1 Representation of Causal Relations

Graphical causal models intend to describe with a Directed Acyclic Graph (DAG) the structure of the underlying physical mechanisms governing a system under study. The state of each variable, represented by a node in the graph, is generated by a stochastic process that is determined by the values of its parent variables in the graph. All variables that influence the outcome of the process are called *causes* of the outcome variable. An *indirect cause* produces the state of the effect indirectly, through another variable. If there is no intermediate variable among the known variables, the cause is said to be a *direct cause*.

Each process represents a physical mechanism. In its most general form it can be described by a conditional probability distribution (CPD)  $P(X \mid Pa(X))$ , where  $Pa(X)$  is the set of parent nodes of  $X$  in the graph and constitute the direct causes of the variable. A causal model consists of a DAG over all variables and a CPD for each variable. The combination of the CPDs results in a joint probability distribution:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid Pa(X_i)) \quad (1.6)$$

For a discrete variable, the CPD is encoded by means of a tabular representation: for each possible assignment of values to the parents of  $X_i$ , we need to specify a distribution over the values that  $X_i$  can take. This is called a conditional probability table. For continuous variables, one often relies on prior knowledge or assumptions about the structure of the distribution. If one assumes linearly-related variables, the CPDs can be described by the following *structural equations*:

$$P(X_i | Pa(X_i)) = \sum_{X_j \in Pa(X_i)} a_{i,j} \cdot X_j + U_i + c_i \quad (1.7)$$

where  $U_i$  represent the stochastic variations which cannot be explained by the model and  $c_i$  a constant term. One often assumes that  $U_i$  is normally distributed.

### 1.3.2 Modularity and the Effect of Changes to the System

A causal model represents a collection of processes that could account for the generation of the observed data. Each process is a stable and autonomous physical mechanism. It is then conceivable to change one such relationship without changing the others. This *modularity* permits one to predict the effect of external interventions or local reconfigurations of the mechanisms (Pearl, 2000). An *intervention* is defined as an atomic operation that fixates a variable to a given state and eliminates the corresponding factor (CPD) from the factorization (Eq. 1.6) (Pearl, 2000). Applied on a causal graph, an intervention on variable  $X$  sets the value of  $X$  and breaks all of the edges in the graph directed into  $X$  and preserves all other edges in the graph, including all edges directed out of  $X$ . This is called the Manipulation Theorem (Spirtes et al., 1993, p. 51). Intervening on a variable only affects its effects. Causes have to be regarded as if they were levers which can be used to manipulate their effects.

This approach does not directly define causality, but defines the implications of having a thorough knowledge of the mechanisms that make up a system. Manipulability puts a constraint of independentness on the mechanisms. The accuracy of the mutilated model relies on autonomy or modularity; a mechanism can be replaced by another without affecting the rest of the system. It is defined by Hausman and Woodward (1999, p. 545) as follows. They relate each CPD to a structural equation (Eq. 1.7).

**Definition 1.1** (*Modularity*) *For all subsets  $\mathbf{Z}$  of the variable set  $\mathbf{V}$ , there is some non-empty range  $\mathbf{R}$  of values of members of  $\mathbf{Z}$  such that if one intervenes and sets the value of the members of  $\mathbf{Z}$  within  $\mathbf{R}$ , then all equations except those with a member of  $\mathbf{Z}$  as a dependent variable (if there is one) remain invariant.*

### 1.3.3 Representation of Independencies

The key for causal inference is the conditional independencies entailed by the system's causal structure. They are based on the property of Markov chains and v-structures. If  $X$  is affected by  $Y$  and  $Z$ , then we do not expect that  $X$  is independent of  $Y$  conditional on  $Z$ , except if  $Y$  affects  $X$  via  $Z$ . This is represented by a Markov chain. Random variables  $X, Z, Y$  are said to form a *Markov chain* in that order, denoted by  $X \rightarrow Z \rightarrow Y$ , if the joint probability mass function can be written as

$$P(X, Z, Y) = P(X).P(Z | X).P(Y | Z) \quad (1.8)$$

which is equivalent to the conditional independence of  $X$  and  $Y$  given  $Z$ . *Conditional independence* of  $X$  and  $Y$  given  $Z$ , written as  $X \perp\!\!\!\perp Y | Z$ , is defined as

$$P(X, Y | Z) = P(X | Z).P(Y | Z) \quad (1.9)$$

The conditional independence expresses that learning the value of  $X$  does not provide additional information about  $Y$  once the state of  $Z$  is known. We say that  $Z$  ‘screens off’  $X$  from  $Y$ . Once the state of  $Z$  is observed, the state of  $Y$  no longer depends on that of  $X$ . For a *v-structure* on the other hand, for example  $X \rightarrow Z \leftarrow Y$ ,  $X$  and  $Y$  are independent, but become dependent when conditioned on  $Z$ .

For a causal model, the *Causal Markov Condition* gives us the independencies that follow from the causal structure: each variable is probabilistically independent of its non-effects conditional on its direct causes. This condition is defined by Spirtes et al. (1993) as follows:

**Definition 1.2** (*Causal Markov Condition*) *Let  $G$  be a causal graph with vertex set  $V$  and  $P$  be a probability distribution over the vertices in  $V$  generated by the causal structure represented by  $G$ .  $G$  and  $P$  satisfy the Causal Markov Condition if and only if for every  $W$  in  $V$ ,  $W$  is independent of  $V \setminus \text{Descendants}(W) \setminus \text{Parents}(W)$  given  $\text{Parents}(W)$ .*

These independencies are irrespective of the nature of the mechanisms, of the exact parameterization of the conditional probability distributions  $P(X_i | Pa(X_i))$ . Pearl and Verma constructed a graphical criterion, called *d-separation*, for retrieving, from the causal graph, all independencies following from the Causal Markov Condition.

A graph is called *faithful* to a distribution if all conditional independencies of the distribution correspond to a *d-separation* in the graph and vice versa. In other words, faithfulness means that if a graph represents a causal structure, all conditional independencies follow from the system’s causal structure.

#### 1.3.4 Correspondence with Bayesian networks

Graphical causal models provide a probabilistic account of causality (Spohn, 2001). This resulted in a close correspondence with Bayesian networks. In contrast to causal models, Bayesian networks are only concerned with offering a dense and manageable representation of joint distributions. A joint distribution over  $n$  variables can be *factorized* relative to a variable ordering  $(X_1, \dots, X_n)$ :

$$P(X_1, \dots, X_n) = \prod_i^n P(X_i | X_1, \dots, X_{i-1}) \quad (1.10)$$

Variable  $X_j$  can be removed from the conditioning set of variable  $X_i$  if it becomes conditionally independent from  $X_i$  by conditioning on the rest of the set:

$$P(X_i | X_1 \dots X_{i-1}) = P(X_i | X_1 \dots X_{j-1}, X_{j+1} \dots X_{i-1}). \quad (1.11)$$

Such conditional independencies reduce the complexity of the factors in the factorization. The conditioning sets of the factors can be described by a Directed Acyclic



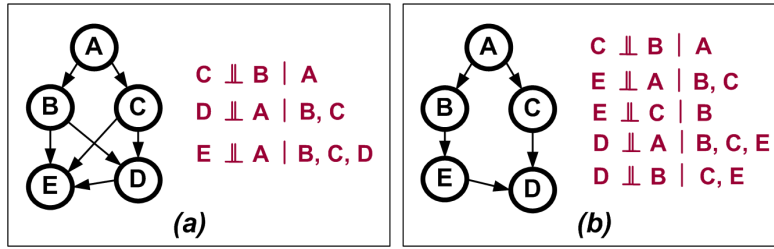


FIG. 1.2. Factorization of the same distribution according to variable ordering ( $A, B, C, D, E$ ) and reduction by three independencies (a), and according to variable ordering ( $A, B, C, E, D$ ) and five independencies (b).

Graph (DAG), in which each node represents a variable and has incoming edges from all variables of the conditioning set of its factor. The joint distribution is then described by the DAG and the conditional probability distributions (CPDs) of the variables conditional on their parents. A *Bayesian network* is a factorization that is edge-minimal, in the sense that no edge can be deleted without destroying the correctness of the factorization.

Although edge-minimality of a Bayesian network, the graph depends on the chosen variable ordering. Some orderings lead to the same networks, while others result in different topologies. Take 5 stochastic variables  $A, B, C, D$  and  $E$ . Figure 1.2(a) shows the graph that was constructed by simplifying the factorization based on variable ordering ( $A, B, C, D, E$ ) by the three given conditional independencies. However, the Bayesian network, describing the same distribution, but based on ordering ( $A, B, C, E, D$ ), depicted in Fig. 1.2(b) contains 2 edges less because of 5 useful independencies. Both networks represent the probabilities just as well, except that the first one is more complex. We call the *minimal factorization* as the factorization which has the least total number of variables in the conditioning sets. The corresponding Bayesian network is called the *minimal Bayesian network* of a probability distribution.

Analogous to the Causal Markov Condition, the Markov Condition gives the conditional independencies that follow from the structure of a Bayesian network: each variable is independent from all its non-descendants by conditioning on its parents in the graph. The equivalence of the Markov Condition and factorizability can be proven (Hausman and Woodward, 1999, p. 532). This ensures the correspondence: causal models are also Bayesian networks. The difference lies in the causal component; causal models attribute a causal interpretation to the edges of the graph and are therefore called *causally interpreted Bayesian networks*.

### 1.3.5 Causal Inference

The goal of causal inference is to learn the causal structure of a system based on observational data. Causal structure learning algorithms fall apart in two categories: scoring-based and constraint-based algorithms.

Scoring-based algorithms are based on an optimized search through the set of all possible models, which tries to find the minimal model that best describes the data. Each model is given a score that is a trade-off between model complexity and goodness-of-fit.

Different scoring criteria have been applied in these algorithms, such as a Bayesian scoring method (Cooper and Herskovits, 1992)(Heckerman et al., 1994), an entropy based method (Herskovits, 1991), a Minimum Message Length (MML) method (Oliver et al., 1992), and one based on the Minimum Description Length (MDL) (Suzuki, 1996). As explained in the introduction, we are not investigating how to select the minimal model from the a priori chosen model class, but the model class which should be considered.

Constraint-based learning algorithms rely on the conditional independencies detected that follow from the system’s causal structure. It is a kind of evidence-based construction, the decisions to include an edge and on the edge’s orientation is based on the presence or absence of certain independencies. The algorithms assume the existence of a faithful graph, i.e. that all independencies follow from the causal structure. They also assume that the correct model is the minimal model. Minimality, faithfulness and the Causal Markov Condition give the 3 assumptions that ensure correct learning (Spirtes et al., 1993). The minimality condition is an edge-minimal condition on the true causal graph.

It must be noted that some algorithms, such as the PC algorithm, also require *causal sufficiency*, i.e. that all common causes should be known: variables that are the direct cause of at least two variables. More sophisticated learning algorithms exist that are capable of detecting latent common causes. For now we will not take the presence of latent variables into consideration and discuss the consequences of this in Section 1.8.

#### 1.4 Related Work

The causal interpretation of a Bayesian network and the validity of faithfulness are often criticized (Freedman and Humphreys, 1999; Cartwright, 2001; Williamson, 2005; Hausman and Woodward, 1999). This paper would like to contribute to the discussion by giving an additional viewpoint through the concept of the KMSS. Some of the examples on which criticism on the possibility of causal inference is based will be discussed in Section 1.8. Hausman and Woodward (1999) on the other hand are strong defenders of linking the causal interpretation of models to modularity. They defend the equivalence of modularity and the Causal Markov Condition (Hausman and Woodward, 1999, p. 554). We will contribute to the discussion by motivating why and when modularity is a valid assumption, and showing the limitations of assuming faithfulness.

Pearl and others use *stability* as the main motivation for the faithfulness of causal models (Pearl, 2000, p. 48). Consider the model of Fig. 1.2(b). In general, one expects  $A$  to depend upon  $D$ .  $A$  and  $D$  are independent only if the stochastic parameterization is such that the influences via paths  $A \rightarrow B \rightarrow E \rightarrow D$  and  $A \rightarrow C \rightarrow D$  cancel out exactly. This system is called unstable because a small change in the parameterization results in a dependency. The unhappy balancing act is a measure zero event, the probability of such a coincidence can therefore be regarded as zero. Hence, the majority of distributions compliant with a DAG are faithful (Pearl, 2000, p. 18). We argue that indeed typical distributions are faithful, but that nonetheless, unfaithful distributions appear.

Milan Studeny was one of the first to point out that the Bayesian networks cannot represent all possible sets of independencies. He constructed a different framework, called

*imsets* (Studeny, 2001), which is capable of representing broader sets of independencies. We advocate a different approach. We will not look for a different representation of conditional independencies, but stick to Bayesian networks. Yet, we will try to find explanations (referring to regularities) for the presence of conditional independencies not coming from the system's causal structure.

### 1.5 Minimal Description of Distributions

In this section we start the analysis of causal inference by applying the KMSS principle on observed data of a collection of independent and identically distributed random variables. A minimal description for the data corresponds to the construction of an efficient code which on its turn corresponds to the description of a probability distribution (Grünwald et al., 2005)[Chapter 2]. We thus have to investigate how distributions can be described compactly.

From the theory of Bayesian networks (Section 1.3.4), we know that a joint distribution can be described shorter by a factorization that is reduced by conditional independencies of the form of Eq. 1.11. The minimal factorization leads to  $P(X_1, \dots, X_n) = \prod CPD_i$ , with  $CPD_i$  the CPD of variable  $X_i$ . The descriptive size of the CPDs is determined by the number of variables in the conditioning sets. A two-part description of a joint distribution is then:

$$\begin{aligned} descr(P(X_1 \dots X_n)) &= descr(\{Pa(X_1), \dots, Pa(X_n)\}) \\ &\quad + descr(CPD_1) + \dots + descr(CPD_n) \end{aligned} \quad (1.12)$$

With  $descr(x)$  the description of  $x$ . Note that the parents' lists are described very compactly by a DAG. If the description according to eq. 1.12 is shorter than the literal description of the joint distribution, then the reduction of the factorization contains meaningful information. This meaningful information is described by the parents' lists or the DAG of the Bayesian network.

**Theorem 1.3** *If the two-part code description of a probability distribution, given by Eq. 1.12 in which the CPDs are described literally, results in an incompressible string which is shorter than the literal description of the joint probability distribution, the first part is the Kolmogorov minimal sufficient statistic.*

**Proof** The CPDs do not contain meaningful information (regularities), since they are literally described and they are incompressible. This last follows from the incompressibility of the total description. Since the total description is shorter than the literal description, the reduction of the factorization outweighs the description of the parents' lists. The parents' lists therefore contain meaningful information. Their incompressibility ensures that it is the KMSS.  $\square$

Concluding, we end up with a three-part code for the description of the observations:

$$\begin{aligned} descr(data) &= descr(DAG) + descr(CPD_1) + \dots + descr(CPD_n) \\ &\quad + descr(data \mid distribution) \end{aligned} \quad (1.13)$$

The data is described with the help of a probability distribution, which on its turn is described by a DAG and a list of CPDs. The regularities that allow the compact description

of the data are the dependencies among the variables; knowing one variable gives information about the state of another variable. Conditional independencies, on the other hand, reduce the model's complexity. They reduce the number of variables to consider when describing the dependencies among the variables.

## 1.6 Minimality of Bayesian Networks

The following two theorems show that the Bayesian network corresponding to the minimal factorization is the KMSS and faithful if its DAG and CPDs are random and incompressible.

**Theorem 1.4** *If a faithful Bayesian network exists for a distribution, it is the minimal Bayesian network, i.e. the Bayesian network with the minimal number of edges.*

**Proof** Recall that the absence of an edge between two variables  $X$  and  $Y$  in a Bayesian network implies that there exists a set of variables  $S$  not containing  $X$  and  $Y$  that makes  $X$  and  $Y$  conditionally independent:  $X \perp\!\!\!\perp Y \mid S$ . In case of faithfulness, the presence of an edge forbids the existence of such a set. Let  $A$  be a graph that has fewer edges than the faithful graph  $B$ . It follows that  $B$  contains an edge between two variables  $X$  and  $Y$  that  $A$  does not contain. The absence of the edge in  $A$  implies that  $X$  and  $Y$  become independent by conditioning on some set of the other variables. But this contradicts with the faithfulness of  $B$  which implies that  $X$  and  $Y$  cannot become independent.  $\square$

Neapolitan (2003)[p. 107] provides a proof for edge-minimality, while here minimality in the global sense is considered.

The DAG of a Bayesian network corresponds to a set of conditional independencies. Intuitively we would expect that two variables are dependent if they are not  $d$ -separated. When this is true, the DAG is faithful to the probability distribution. The next theorem proves that two variables that are not  $d$ -separated can only be independent if there is a constraint between the probabilities. To illustrate the theorem, consider the model of Fig. 1.2(b) and the set of distributions compatible with the DAG. For typical distributions, dependencies  $D \not\perp\!\!\!\perp E$  and  $A \not\perp\!\!\!\perp E$  hold. There are, however, specific parameterizations which lead to independencies  $D \perp\!\!\!\perp E$  or  $A \perp\!\!\!\perp E$ . Such independencies only follow if specific equations between the free parameters are satisfied.

**Theorem 1.5** *A Bayesian network for which the concatenation of the descriptions of the conditional probability distributions (CPDs) is incompressible, is faithful.*

**Proof** Recall that a Bayesian network is a factorization that is edge-minimal. This means that for each parent  $pa_{i,j}$  of variable  $X_i$ :

$$P(X_i \mid pa_{i,1}, \dots, pa_{i,j}, \dots, pa_{i,k}) \neq P(X_i \mid pa_{i,1}, \dots, pa_{i,j-1}, pa_{i,j+1}, \dots, pa_{i,k}) \quad (1.14)$$

Variables cannot be eliminated from the factors of the factorization. The proof will show that any two variables that are not  $d$ -separated are dependent, unless the probabilities of the CPDs are 'related', in the sense that some probabilities can be calculated from others and the set of CPDs is compressible. We derive the relations for discrete variables. For

continuous variables, the analysis results in relations among the free parameters of the CPDs.

We have to consider the following possibilities. The two variables can be adjacent (a), related by a Markov chain (b)<sup>1</sup>, a v-structure (c), a combination of both or connected by multiple paths (d).

First we prove that a variable marginally depends on each of its adjacent variables (a). Consider adjacent nodes  $D$  and  $E$  of the Bayesian network of Fig. 1.2(b). We will demonstrate that  $P(D | E) = P(D)$  results in a regularity. We expand the first term with all other parents of  $D$ :

$$P(D | E) = \sum_{c \in C_{dom}} P(D | E, c).P(c | E) \quad (1.15)$$

$C$  is also a parent of  $D$ , thus, by Eq. 1.14, there are at least two values of  $C_{dom}$  for which  $P(D | E, c) \neq P(D | E)$ <sup>2</sup>. Take  $c_1$  and  $c_2$  being such values for which

$$P(D | E, c_1) \neq P(D | E, c_2). \quad (1.16)$$

There are also at least 2 such values of  $E_{dom}$ , take  $e_1$  and  $e_2$ . Eq. 1.15 should hold for all values of  $E$  and equal to  $P(D)$  to get an independency. This results in the following relation among the probabilities:

$$\begin{aligned} & P(D | e_1, c_1).P(c_1 | e_1) + P(D | e_1, c_2).P(c_2 | e_1) \\ &= P(D | e_2, c_1).P(c_1 | e_1) + P(D | e_2, c_2).P(c_2 | e_1) \end{aligned} \quad (1.17)$$

Note that the equation cannot be algebraically simplified: the conditional probabilities are not equal to  $P(D)$  (Eq. 1.14) nor to each other (Eq. 1.16). The proof can easily be generalized for variables having more parents.

Next, by the same arguments it can be proved that variables connected by a Markov chain are by default dependent (b). Take  $A \rightarrow B \rightarrow E$  in Fig. 1.2(b), independence of  $A$  and  $E$  requires that

$$P(E | a) = \sum_{b \in B_{dom}} P(E | b).P(b | a) = P(E) \quad \forall a \in A. \quad (1.18)$$

and this would also result in a regularity among the CPDs.

In a v-structure, both causes are dependent when conditioned on their common effect (c), for  $C \rightarrow D \leftarrow E$ ,  $P(D | C, E) \neq P(D | E)$  is true by Eq. 1.14. Finally, if there are multiple unblocked paths connecting two variables, then independence of both variables implies a regularity as well (d). Take  $A$  and  $D$  in Fig. 1.2(b):

$$P(D | A) = \sum_{b \in B_{dom}} \sum_{c \in C_{dom}} \sum_{e \in E_{dom}} P(D | c, e).P(c | A).P(e | b).P(b | A).$$

<sup>1</sup>Recall that a Markov chain is a path not containing v-structures.

<sup>2</sup> $P(D | E)$  is a weighted average of  $P(D | E, C)$ . If one probability  $P(D | E, c_1)$  is different than this average, let's say higher, than there must be at least one value lower than the average, thus different.

Note that  $P(c, e | A) = P(c | A) \cdot P(e | A)$  follows from the independence of  $C$  and  $E$  given  $A$ . All factors from the equation satisfy Eq. 1.14, so that, again, the equation only would equal to  $P(D)$  if there is a relation among the probabilities.  $\square$

From the theorem it follows that a Bayesian network with random CPDs is the minimal factorization. Bayesian networks not based on a minimal factorization, such as the one of Fig. 1.2, are compressible, namely by the regularities among the CPDs that follow from the independencies not represented by the graph. Pearl hypothesizes that there is no bounded set of conditions that would ensure the existence of a faithful graph (Pearl, 1988, p. 131). Indeed, as shown by the theorem, every dependence can be turned into an independence by a balanced parameterization of some CPDs.

It must be noted that if there exists a faithful Bayesian network, it is not necessarily unique. Multiple faithful models can exist for a distribution. These models represent the same set of independencies and are therefore statistically indistinguishable. They define a *Markov-equivalence class*. It is proved that they share the same skeleton and v-structures. They only differ in the orientation of some edges (Pearl, 2000). This set can be represented by a partially-directed acyclic graph in which some of the edges are not oriented. The corresponding factorizations have the same number of conditioning variables and thus all models of a Markov-equivalence class have the same complexity.

### 1.7 When the Minimal Bayesian Network is the KMSS

In this section we will discuss the case in which there is exactly one minimal Bayesian network which is also the minimal description. This means that there are no other regularities and no other independencies than the conditional independencies represented by the model. The DAG is then the KMSS and minimally represents all regularities. It is also faithful.

The minimal Bayesian network decomposes the description of a joint distribution into a list of CPDs. This means that the minimal description of the system is a concatenation of descriptions, namely the description of the individual CPDs. In other words, we have found a unique and minimal decomposition of the model. This brings us to modularity and manipulability. We have discovered that the minimal description is a concatenation of unrelated components. The CPDs are independent; the concatenation of their descriptions cannot be compressed. Then, among all possible explanations, the simplest is that *each CPD corresponds to an independent part of reality*. Thus, following Occam's Razor, modularity is the most likely hypothesis about the system under study. The correctness of Occam's razor cannot be proven, the principle must be interpreted as the most effective *strategy* for deciding among competing explanations (Grünwald, 1998). Modularity of the minimal Bayesian network must be regarded as the top-ranked hypothesis, which can be verified with background knowledge or experiments with interventions. Thus, the three conditions for causal inference are valid (Section 1.3.5): minimality and faithfulness are fulfilled, and the Causal Markov Condition follows from modularity. Description minimality is linked to causality through modularity.

Occam's razor is contradicted when the real system is more complex than suggested by the complexity of the observations. Take the impact of *Tax rate* increase on *Tax revenue* as shown in Fig. 1.3(a). A *Tax rate* increase has a negative effect on the

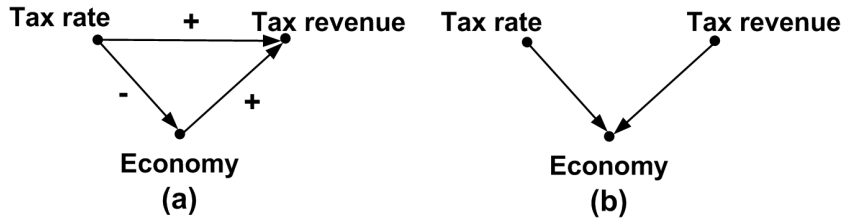


FIG. 1.3. Model in which the impact of *Tax rate* increase on *Tax revenue* is neutralized by the negative effect on the *Economy* (a). The minimal Bayesian network describing the system(b).

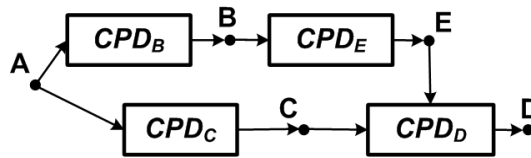


FIG. 1.4. Decomposition of the system represented by the causal model of Fig. 1.2(b) into independent components.

*Economy* which could neutralize the increase of the tax revenues, such that *Tax rate*  $\perp\!\!\!\perp$  *Tax revenue*. If so, the system is minimally described by the model of Fig. 1.3(b). This model is faithful, incompressible and simpler than the true model. From observations alone, one cannot find indications for the more complex true model. Although not minimal in the global sense, the model of Fig. 1.3(a) is edge-minimal: no edge can be removed without destroying the correctness of the model.

The CPD of a variable is also called the variable's *causal Markov kernel*. Note that by representing a causal model with a graph, the representation suggests that the edges - instead of the CPDs - are the basic components. This is however not true. A graphical model can therefore be misleading. A better representation is shown in Fig. 1.4. It represents the same system as the causal model of Fig. 1.2(b), but emphasizes that CPDs are the basic components.

Decomposition and thus also causality matches with a *reductionist* view, according to which the world can be studied in parts. Indeed, if the system cannot be decomposed, if there are no conditional independencies that simplify a factorization, then the DAG does not contain meaningful information. We end up with a *Holist* system in which everything depends on everything.

Note that uniqueness of the minimal Bayesian network is not essential. As discussed in the previous section, if the minimal Bayesian network is not unique, the Markov-equivalence class indicates exactly which parts are undecided (the orientation of some edges). So, we know exactly for which parts of the model we do not have enough information to decide upon the decomposition.

### 1.8 When the Minimal Bayesian Network is not the KMSS

To study the validity of faithfulness and the modularity property, we will in this section not assume incompressibility of the minimal Bayesian networks. They are denoted with  $BN_{min}$ . Instead, we will study a wide variety of cases, appearing throughout literature, in which regularities appear that are not described by a Bayesian network. We will analyze the properties of the True Causal Model ( $CM_{true}$ ) and those of the  $BN_{min}$ .

Table 1 gives an overview of the answers for the next questions, which will be discussed in the following.

- Is the  $CM_{true}$  compressible? If so, is the compressibility due to the compressibility of the description of a single CPD or the compressibility of the concatenation of the descriptions of multiple CPDs?
- Is the compressibility of the minimal Bayesian networks due to the compressibility of the description of a single CPD or the compressibility of the concatenation of the descriptions of multiple CPDs?
- Is the  $CM_{true}$  present in  $BN_{min}$ ? The answer to this and the next question determines the feasibility of causal inference.
- Is there a unique  $BN_{min}$ ? Are the regularities under consideration responsible for the presence of multiple minimal Bayesian networks?
- Is the true causal model faithful to the system?
- Are the minimal Bayesian networks faithful to the system?
- Does modularity holds for the true causal model?

	Compress. $CM_{true}$	Compress. $BN_{min}$	$CM_{true}$ $\in BN_{min}$	Unique $BN_{min}$	$CM_{true}$ faithful	$BN_{min}$ faithful	Modular $CM_{true}$
1. Local	single	single	Yes	Yes	Yes	Yes	Yes
2. PIM	single	single	Yes	No	No	No	Yes
3. Determ	single	single	Yes	No	No	No	Yes
4. Unfaithf	concat.	concat.	Yes	Yes	No	No	No/Yes
5. Markov	No'	concat.	No	No	Yes'	No	Yes'
6. Latent	No*	concat.	No	No	Yes*	No	Yes*
7. OO-nets	concat.	concat.	Yes	Yes	Yes	Yes	Yes

TABLE 1.1. Answers to questions for the different case studies. A ' indicates that Markov networks are considered. An asterisk (\*) indicates that Bayesian networks with latent variables are considered.

#### 1.8.1 Compressibility of a single CPD

First we consider cases in which the description of an individual CPD is compressible. Faithfulness and the uniqueness of the minimal Bayesian network are not guaranteed, but the cases show that the modularity assumption still holds. The CPDs are independent.



**Case 1.** When individual CPDs can be compressed, we call this type of regularity *local structure* (Friedman and Goldszmidt, 1996). For discrete variables, the conditional probability tables are exponential in the number of parents of a variable  $X$ : for each possible assignment of values to the parents of  $X$ , we need to specify a distribution over the values  $X$  can take. When regularities among the probabilities appear these tables can be described more compactly, for example by decision trees. The regularities to construct the tree are called context-specific independencies (Boutilier et al., 1996). On top of the independencies following from the causal structure, the system exhibits additional regularities. But the model remains faithful and the decomposition is correct.

**Case 2.** Variables in *pseudo-independent models* are pairwise independent but collectively dependent (Xiang et al., 1996). For example, consider a binary variable  $X_3$  that is determined by two other binary variables  $X_1$  and  $X_2$  by an *exclusive or* relation:  $X_3 = X_1 \text{ EXOR } X_2$ . This system can be represented by causal model  $X_1 \rightarrow X_3 \leftarrow X_2$ . Because of the pairwise independencies  $X_3 \perp\!\!\!\perp X_1$  and  $X_3 \perp\!\!\!\perp X_2$ , the model is not faithful. There are three minimal Bayesian networks: besides the correct  $X_1 \rightarrow X_3 \leftarrow X_2$ , also  $X_1 \rightarrow X_2 \leftarrow X_3$  and  $X_2 \rightarrow X_1 \leftarrow X_3$ . The CPD  $P(X_3 | X_1, X_2)$  exhibits a strict regularity. Yet, pseudo-independent models fit in the reductionist approach of causal models. The only problem is that the conditional independencies do not provide enough information to conclude about the causal connections.

**Case 3.** Deterministic or functional relations among variables result in CPDs with a very specific form. Distributions with deterministic relations cannot be represented by a faithful graph (Spirtes et al., 1993). Consider the system  $X \rightarrow Y \rightarrow Z$  in which  $Y$  is a function of  $X$ :  $Y = f(X)$ . From the model (Markov chain) it follows that  $X \perp\!\!\!\perp Z | Y$ . By the functional relation, variable  $X$  got all information about  $Y$ , which implies  $Y \perp\!\!\!\perp Z | X$ . Both independencies imply a violation of the *intersection condition*, one of the conditions that Pearl imposes on a distribution in the elaboration of causal theory and its algorithms (Pearl, 1988). In (Lemeire, 2007) we call  $X$  and  $Y$  *information equivalent* with respect to  $Z$ , both variables have in some sense the same information about  $Z$ . Then, the set of minimal Bayesian networks contains graphs that connect  $X$  with  $Z$  and graphs that connect  $Y$  with  $Z$ . From the information about the conditional independencies alone we cannot decide upon which variable,  $X$  or  $Y$ , directly relates to  $Z$ . The solution we proposed for causal inference is to connect the variables that have the simplest relation (Lemeire, 2007). We defined an augmented causal model which also incorporates information of deterministic relations.

### 1.8.2 Compressibility of a set of CPDs

When the description of some CPDs taken together can be compressed, the CPDs are in some way related.

**Case 4.** The most-known example of unfaithfulness is when in the model of Fig. 1.5 (a),  $A$  and  $D$  appear to be independent (Spirtes et al., 1993). This happens when the influences along the paths  $A \rightarrow B \rightarrow D$  and  $A \rightarrow C \rightarrow D$  exactly balance, so that they cancel each other out and the net effect results in an independence. For continuous variables this happens when an exact correspondence of the free parameters is fulfilled. The model is not faithful. This balancing act can give an indication of a global mecha-

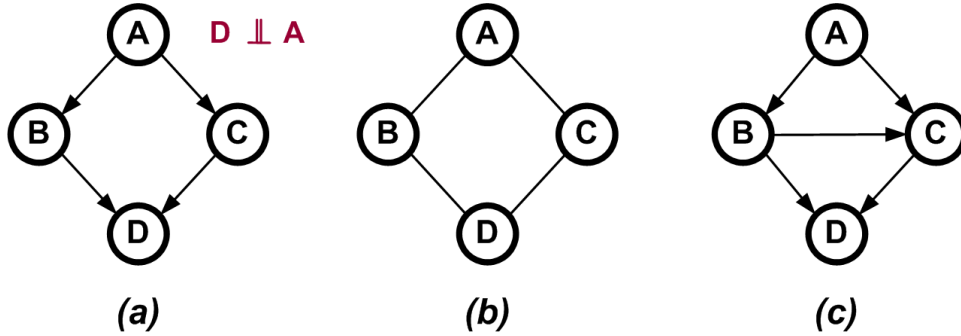


FIG. 1.5. O-structure in which  $A$  is independent from  $D$  (a). A Markov network (b) and one of the minimal Bayesian networks describing the same system (c).

nism or *meta-mechanism*, such as evolution (Korb and Nyberg, 2006), controlling the mechanisms such that the parameters are calibrated until they neutralize. Modularity and autonomy of the CPDs depends on the meta-mechanism. Evolution works on the long-term, so modularity holds for a limited time period. For meta-mechanisms controlling the mechanisms instantly, the CPDs cannot be considered as being independent.

**Case 5.** Consider a system that is minimally described by a Markov network, as shown in Figure 1.5 (b). Variables which are connected by a path in the network are dependent, unless each path is blocked by the conditioning variables. So is  $B \not\perp C \mid A$ , but  $B \perp C \mid \{A, D\}$ . For describing the same network with a DAG, we have to orient the edges of the network. For acyclicity, we have to create at least one v-structure. We can choose for example  $B - D - C$ . But then, for keeping the same dependencies, we have to add an edge, as shown in Figure 1.5 (c). Without  $B \rightarrow C$  we would have  $B \perp C \mid A$ . Clearly, this Bayesian network is not minimal; the description is longer than that of a Markov network. The parameterizations of the CPDs contain redundancies. In the model of 1.5 (c), the parameterizations must ensure that  $B \perp C \mid \{A, D\}$ , an independency which is not captured by the DAG. The causal interpretation of the CPDs (modularity) is not correct for the minimal Bayesian networks.

**Case 6.** Causal sufficiency, the knowledge of all common causes, is an important property for correct causal learning. Take the system depicted in Fig. 1.6(a) in which  $L$  is an unknown variable which is the cause of  $B$  and  $C$ . This gives rise to multiple minimal Bayesian networks, none of which models the system correctly. One of them is depicted in Fig. 1.6(b).  $B$  and  $C$  are correlated, but none of the other known variables is the cause of both, so either  $B$  should be oriented towards  $C$  or vice versa.  $A$  should be connected to  $C$  to reflect dependency  $A \not\perp C \mid B$ . But  $A \perp C$ , thus there is a dependency between  $P(B \mid A)$  and  $P(C \mid A, B, D)$ . The Bayesian network is therefore compressible and not faithful ( $A \perp C$  is not represented). The solution is to look for an alternative model class. Spirtes et al. (1995) propose the use of a *Partially-oriented Acyclic Graph* (PAG) by which one can express the possibility of latent variables.

**Case 7.** Another regularity is the repetition of similar mechanisms in a system. This results in a causal model in which identical CPDs appear. The model is therefore com-

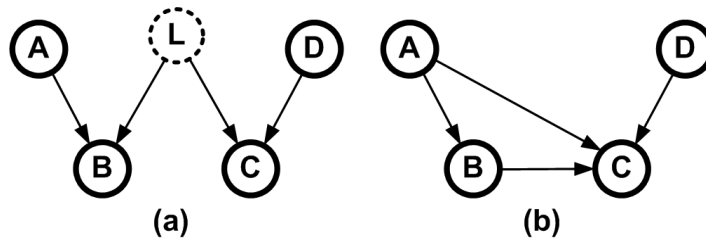


FIG. 1.6. Learning the model of a system with a latent variable  $L$  (a) and one of the minimal Bayesian networks (b).

pressible. The compressibility does not, however, result in a dependence of the CPDs in terms of manipulability. One mechanism can still be replaced by another without affecting the rest of the model. Modularity still holds. *Object-Oriented nets* provide a representation format that explicitly capture similarities of mechanisms (Koller and Pfeffer, 1997).

### 1.9 Conclusions

A Bayesian network decomposes the description of a joint probability distribution into Conditional Probability Distributions (CPDs). If a Bayesian network provides the Kolmogorov Minimal Sufficient Statistic (KMSS) of a system, it gives the most plausible hypothesis about the causal structure of the system. The CPDs can be matched up with mechanisms of the underlying system. Decomposition reflects the causal component of graphical causal models.

Causal model theory expresses what typically can be expected from a causal structure. Typical distributions that are compatible with a causal structure are faithful. However, atypical distributions contain additional regularities and may invalidate the above conclusions. The minimal Bayesian networks of a probability distribution are then compressible and do not represent the KMSS.

If the description of a single CPD is compressible, this can result in unfaithfulness of the causal model. Causal inference is still possible, since the true model is an element of the set of minimal Bayesian networks and modularity is plausible. If on the other hand the concatenation of the CPDs is compressible, then the CPDs are no longer independent and the mapping of CPDs onto independent mechanisms becomes invalid. This can be due to a kind of meta-mechanism governing other mechanisms, or the incorrectness of considering the set of Bayesian networks as model class.

## Bibliography

- Craig Boutilier, Nir Friedman, Moises Goldszmidt, and Daphne Koller. Context-specific independence in Bayesian networks. In *Uncertainty in Artificial Intelligence*, pages 115–123, 1996.
- Nancy Cartwright. What is wrong with Bayes nets? *The Monist*, pages 242–264, 2001.
- G.F. Cooper and E. Herskovits. A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, 9:309–347, 1992.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, Inc., 1991.
- D.A. Freedman and P. Humphreys. Are there algorithms that discover causal structure? *Synthese*, 121:2954, 1999.
- N. Friedman and M. Goldszmidt. Learning bayesian networks with local structure. In *In Proc. of the 12th Conf. on Uncertainty in Artificial Intelligence*, 1996.
- Péter Gács, J. Tromp, and Paul M. B. Vitányi. Algorithmic statistics. *IEEE Trans. Inform. Theory*, 47(6):2443–2463, 2001.
- P. Grünwald. *The Minimum Description Length Principle and Reasoning under Uncertainty*, *ILLC Dissertation series 1998-03*. PhD thesis, University of Amsterdam, 1998.
- P. Grünwald, I.J. Myung, and M.A. Pitt. *Advances in Minimum Description Length Principle. Theory and Applications*. MIT Press, 2005.
- Daniel M. Hausman and James Woodward. Independence, invariance and the causal Markov condition. *British Journal For the Philosophy Of Science*, 50(4):521–583, 1999.
- D. Heckerman, D. Geiger, and D.M. Chickering. Learning Bayesian networks: the combination of knowledge and statistical data. Technical Report MSR-TR-94-09, Microsoft Research, 1994.
- E.H. Herskovits. *Computer-Based Probabilistic Network Construction*. PhD thesis, Medical information sciences, Stanford University, CA, 1991.
- Daphne Koller and Avi Pfeffer. Object-oriented Bayesian networks. In *Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence (UAI-97)*, pages 302–313, 1997.
- Kevin B. Korb and Erik Nyberg. The power of intervention. *Minds and Machines*, 16(3):289–302, 2006.
- Jan Lemeire. *Learning Causal Models of Multivariate Systems and the Value of it for the Performance Modeling of Computer Programs*. PhD thesis, Vrije Universiteit Brussel, 2007.
- Ming Li and Paul M. B. Vitányi. *An Introduction to Kolmogorov Complexity and Its Applications*. Springer Verlag, 1997.
- Richard E. Neapolitan. *Learning Bayesian Networks*. Prentice Hall, 2003.

- J. J. Oliver, D. L. Dowe, and C. S. Wallace. Inferring decision graphs using the minimum message length principle. In *Proc. fifth Australian joint conf. on artificial intelligence, Tasmania, Australia*, pages 361–367, 1992.
- Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Mateo, CA, Morgan Kaufman Publishers, 1988.
- Judea Pearl. *Causality. Models, Reasoning, and Inference*. Cambridge University Press, 2000.
- P. Spirtes, C. Meek, and T. Richardson. Causal inference in the presence of latent variables and selection bias. In *Proc. of the 11th Conference on Uncertainty in Artificial Intelligence, ed P. Besnard and S. Hanks*, pages 499–506. Morgan Kaufmann, 1995.
- Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search*. Springer Verlag, 2nd edition, 1993.
- Wolfgang Spohn. Bayesian nets are all there is to causal dependence. In *In Stochastic Causality, Maria Carla Galavotti, Eds. CSLI Lecture Notes*, 2001.
- Milan Studeny. On non-graphical description of models of conditional independence structure. In *HSSS Workshop on Stochastic Systems for Individual Behaviours*, Louvain la Neuve, Belgium, January 2001.
- J. Suzuki. Learning Bayesian belief networks based on the MDL principle: An efficient algorithm using the branch and bound technique. In *Procs of the International Conf. on Machine Learning*, Bally, Italy, 1996.
- Paul M. B. Vitányi. Meaningful information. In Prosenjit Bose and Pat Morin, editors, *ISAAC*, volume 2518 of *Lecture Notes in Computer Science*, pages 588–599. Springer, 2002.
- Jon Williamson. *Bayesian Nets And Causality: Philosophical And Computational Foundations*. Oxford University Press, 2005.
- Yang Xiang, S. K. Wong, and N. Cercone. Critical remarks on single link search in learning belief networks. In *Proceedings of the 12th Annual Conference on Uncertainty in Artificial Intelligence (UAI-96)*, pages 564–571, San Francisco, CA, 1996. Morgan Kaufmann Publishers.