

Markov Blanket and Markov Boundary of Multiple Variables

Xu-Qing Liu

LIUXUQING688@163.COM

*State Key Laboratory of Mechanics and Control of Mechanical Structures
Institute of Nano Science and Department of Mathematics
Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China
Faculty of Mathematics and Physics
Huaiyin Institute of Technology, Huai'an 223003, China*

Xin-Sheng Liu*

XSLIU@NUAA.EDU.CN

*State Key Laboratory of Mechanics and Control of Mechanical Structures
Institute of Nano Science and Department of Mathematics
Nanjing University of Aeronautics and Astronautics, Nanjing 210016, China*

Editors: Marina Meila; Kevin Murphy; Joris Mooij

Abstract

Markov blanket (Mb) and Markov boundary (MB) are two key concepts in Bayesian networks (BNs). In this paper, we study the problem of Mb and MB for multiple variables. First, we show that Mb possesses the additivity property under the local intersection assumption, that is, an Mb of multiple targets can be constructed by simply taking the union of Mbs of the individual targets and removing the targets themselves. MB is also proven to have additivity under the local intersection assumption. Second, we analyze the cases of violating additivity of Mb and MB and then put forward the notions of Markov blanket supplementary (MbS) and Markov boundary supplementary (MBS). The properties of MbS and MBS are studied in detail. Third, we build two MB discovery algorithms and prove their correctness under the local composition assumption. We also discuss the ways of practically doing conditional independence tests and analyze the complexities of the algorithms. Finally, we make a benchmarking study based on six synthetic BNs and then apply MB discovery to multi-class prediction based on a real data set. The experimental results reveal our algorithms have higher accuracies and lower complexities than existing algorithms.

Keywords: Markov blanket, Markov boundary, Markov blanket supplementary, Markov boundary supplementary, Bayesian network

1. Introduction

Bayesian networks (BNs) are graphical structures used to represent the probabilistic relations among a large number of variables and to make the associated probabilistic inferences (Neapolitan, 2004; Pearl, 1988). In recent years, BNs have become one of the most powerful tools in encoding uncertain expert knowledge in expert systems (Daly et al., 2011; Parviainen and Koivisto, 2013) and also deeply influenced on many other actual domains such as medical diagnosis, financial analysis, bioinformatics, and industrial applications (Zhang and Guo, 2006).

*. Corresponding Author.

As two important concepts in BNs, Markov blanket (Mb) and Markov boundary (MB) play a key role in feature selection (FS; Fu and Desmarais, 2010; Pellet and Elisseeff, 2008; Aliferis et al., 2010a,b). Mathematically, Pearl (1988, pp. 218–221) showed the conditional probability for the target given other variables can be replaced by the MB as the conditional set. Pellet and Elisseeff (2008, pp. 1299, 1302) proved that an MB is the theoretically optimal set of features. Further, under certain assumptions about the learner and the loss function, MB is the solution to the variable selection problem (Tsamardinos and Aliferis, 2003; Statnikov et al., 2013).

So far most authors have focused on the problem of Mb or MB for a single variable. In this paper, we consider the problem of Mb and MB for multiple variables. This occurs if, for example, one wants to compute the joint probability of two or more variables conditioned on all other variables. The basic question for Mb of multiple variables is whether the additivity property holds, that is, can an Mb of multiple variables be constructed by simply taking the union of the Mbs of the individual variables and removing the target variables themselves? The same question is for MB. Further, if the additivity property is violated in some situation, how can we do it?

In the literature, there have been lots of MB discovery algorithms, such as the Koller-Sahami (KS) algorithm (Koller and Sahami, 1996), the grow-shrink (GS) algorithm (Margaritis and Thrun, 1999, 2000), the incremental association Markov boundary (IAMB) algorithm (Tsamardinos et al., 2003) and its several variants, the HITON algorithm (Aliferis et al., 2003), the max-min Markov boundary (MMMB) algorithm (Tsamardinos et al., 2006), the parents and children based Markov boundary (PCMB) algorithm and KIAMB algorithms (Peña et al., 2007), the BFMB algorithm (Fu and Desmarais, 2007), the algorithmic framework called generalized local learning (GLL, Aliferis et al., 2010a), and some others (Fu and Desmarais, 2010; Schlüter, 2014). For a single target variable, most of these algorithms are efficient to seek an approximate MB; for multiple target variables, if simply regarding them as a multivariate variable, these algorithms seem to be feasible. However, this will lead to low accuracies and high computational complexities. Hence, it is necessary to design more efficient MB discovery algorithms for multiple variables.

The remainder of this paper is organized as follows. Section 2 presents necessary preliminaries and the motivations of this paper. Subsection 3.1 shows additivity of Mb and MB under the local intersection assumption. In Subsection 3.2, we first analyze when additivity is violated and then put forward the notions of Markov blanket supplementary (MbS) and Markov boundary supplementary (MBS). The properties of MbS and MBS are studied detailedly. In Section 4, we design two MB discovery algorithms for multiple variables, and prove their correctness under the local composition assumption. In addition, we discuss the ways of practically doing conditional independence (CI) tests and analyze the complexities of the algorithms. Section 5 makes a benchmarking study based on six synthetic BNs, and Section 6 considers a practical application. The experimental results show the superiority of our algorithms with higher accuracies and lower complexities than existing algorithms. Section 7 concludes this paper and presents three remarks.

2. Preliminaries and Motivations

In the paper, we denote a variable and its value by upper-case and lower-case letters in italics (e.g., X , x), a set of variables and its value by upper-case and lower-case bold letters in italics (e.g., \mathbf{X} , \mathbf{x}). The difference between \mathbf{X} and \mathbf{Y} is denoted by $\mathbf{X} \setminus \mathbf{Y}$. For brevity, we write $(\mathbf{X} \setminus \mathbf{Y}) \setminus \mathbf{Z}$ as $\mathbf{X} \setminus \mathbf{Y} \setminus \mathbf{Z}$. In addition, we use $|\mathbf{X}|$ to denote the number of variables involved in \mathbf{X} .

2.1 Preliminaries

Suppose we have a joint probability distribution \mathbb{P} over $V \triangleq \{X_1, \dots, X_p\}$ and a directed acyclic graph (DAG) \mathbb{G} with the variables in V as its nodes. We say (\mathbb{G}, \mathbb{P}) satisfies the Markov condition if every $X \in V$ is conditionally independent of its nondescendants given its parents; Further, (\mathbb{G}, \mathbb{P}) is called a Bayesian network (BN) if it satisfies the Markov condition; Furthermore, (\mathbb{G}, \mathbb{P}) satisfies the faithfulness condition if, based on the Markov condition, \mathbb{G} entails all and only conditional independences (CIs) in \mathbb{P} (Pearl, 1988; Neapolitan, 2004).

We write $X \perp\!\!\!\perp Y | Z$ ($X \not\perp\!\!\!\perp Y | Z$), if X and Y are conditionally independent (dependent) given Z with respect to \mathbb{P} . The following properties describe the relations among CI statements (Pearl, 1988; Peña et al., 2007; Statnikov et al., 2013). For any $X, Y, Z, W \subseteq V$, we have (i) *symmetry*: $X \perp\!\!\!\perp Y | Z$ is equivalent to $Y \perp\!\!\!\perp X | Z$; (ii) *decomposition*: $X \perp\!\!\!\perp Y \cup W | Z$ implies $X \perp\!\!\!\perp Y | Z$ and $X \perp\!\!\!\perp W | Z$; (iii) *weak union*: $X \perp\!\!\!\perp Y \cup W | Z$ implies $X \perp\!\!\!\perp Y | Z \cup W$; (iv) *contraction*: $X \perp\!\!\!\perp Y | Z \cup W$ and $X \perp\!\!\!\perp W | Z$ imply $X \perp\!\!\!\perp Y \cup W | Z$; (v) *self-conditioning*: $X \perp\!\!\!\perp Y | Y \cup Z$. Further, if \mathbb{P} is strictly positive, then besides (i)~(v) we also have (vi) *intersection*: $X \perp\!\!\!\perp Y | Z \cup W$ and $X \perp\!\!\!\perp W | Z \cup Y$ imply $X \perp\!\!\!\perp Y \cup W | Z$. Furthermore, if \mathbb{P} is faithful to a DAG \mathbb{G} , then besides (i)~(vi) we also have (vii) *composition*: $X \perp\!\!\!\perp Y | Z$ and $X \perp\!\!\!\perp W | Z$ imply $X \perp\!\!\!\perp Y \cup W | Z$.

Among these properties, intersection and composition are two global ones. Statnikov et al. (2013, p. 504) provided a relaxed version for composition called *local composition*: one says $T \subseteq V$ satisfies the local composition property, if $T \perp\!\!\!\perp X | Z$ and $T \perp\!\!\!\perp Y | Z$ imply $T \perp\!\!\!\perp X \cup Y | Z$ for any $X, Y, Z \subseteq V \setminus T$. We will provide a relaxed version for the intersection property.

Conditional mutual information (CMI) is one of the basic tools for testing CIs. Denote the CMI between X and Y conditioned on Z by $\mathbb{I}(X; Y | Z)$. Then $\mathbb{I}(X; Y | Z) \geq 0$, with equality holding if and only if $X \perp\!\!\!\perp Y | Z$ (Zhang and Guo, 2006). For a practical problem, we cannot access to the true CMI; instead, we use its empirical estimate, denoted by $\mathbb{I}_D(X; Y | Z)$, based on the data D (Cheng et al., 2002). Note that $\mathbb{I}_D(X; Y | Z) \geq 0$ also holds for any $X, Y, Z \subseteq V$.

The *chain rule* for CMI (Cover and Thomas, 2006) is useful to prove the main results of this paper: $\mathbb{I}(X; Y_1 \cup Y_2 | Z) = \mathbb{I}(X; Y_1 | Z) + \mathbb{I}(X; Y_2 | Z \cup Y_1)$ holds for any four sets of variables X, Y_1, Y_2 , and Z from V .

Another notion closely related to CI is d-separation (Pearl, 1988, p. 117). For a DAG \mathbb{G} over V , letting $X, Y, Z \subseteq V$ be disjoint, we say Z d-separates X and Y if it blocks every path between X and Y , and if this is the case we write $X \perp\!\!\!\perp Y | Z$. Here, Z blocking a path c means that c has a head-to-tail node or a tail-to-tail node belonging to Z , or that c has a head-to-head node C such that C and its all descendants are not in Z . As well known, $X \perp\!\!\!\perp Y | Z \Rightarrow X \perp\!\!\!\perp Y | Z$, if (\mathbb{G}, \mathbb{P}) is a BN (Neapolitan, 2004, p. 74). This implication provides a convenient way of identifying CIs.

For example, consider a BN with the graph presented in Figure 1 as its DAG. It follows that: X_2 and X_8 are d-separated by $\{X_4, X_5\}$, meaning $X_2 \perp\!\!\!\perp X_8 | \{X_4, X_5\}$ and thus $X_2 \perp\!\!\!\perp X_8 | \{X_4, X_5\}$; X_3 and X_4 are d-separated by \emptyset , meaning $X_3 \perp\!\!\!\perp X_4$, so $X_3 \perp\!\!\!\perp X_4$. Note that these two probabilistic CIs can not be directly derived from the Markov condition.

In what follows, the concepts of Mb and MB are presented. They are a direct extension of Mb and MB for a single target variable (Pearl, 1988, p. 97; Neapolitan, 2004, pp. 108–109): an Mb of T is a set of variables shielding T from all other variables, so it carries all information of T that cannot be obtained from other variables, while an MB is a minimal Mb.

Definition 1 Let $T \subseteq V$ and $M \subseteq V \setminus T$. We call M a Markov blanket (Mb) of T if $T \perp\!\!\!\perp V \setminus M \setminus T | M$. Further, a Markov boundary (MB) of T is any Mb such that none of its proper subsets is an Mb. ■

When $|T| = 1$, the following results are well known in the literature (Pearl, 1988; Neapolitan, 2004; Statnikov et al., 2013): (a) if (\mathbb{G}, \mathbb{P}) is a BN, then for $T \in V$ the set of its all parents, children, and spouses is an Mb of T (denoted by M_T); (b) if \mathbb{P} satisfies the intersection property, then T has a unique MB; (c) if (\mathbb{G}, \mathbb{P}) satisfies the faithfulness condition, then M_T is the unique MB of T .

Consider again the BN with the graph presented in Figure 1 as its DAG. In this BN, it is seen that $M_{X_4} \triangleq \{X_2, X_6, X_3\}$ is an Mb of X_4 ; further, M_{X_4} is the unique MB of X_4 if the faithfulness condition is satisfied. Similarly, $M_{X_2} \triangleq \{X_4, X_5\}$ is the unique MB of X_2 under the faithfulness condition.

The above result (b) points out that if the uniqueness of MB is violated, then the intersection property must be violated. Lemeire (2007) provided a case of violating intersection called *information equivalence*: X and Y are called information equivalent with respect to T if $T \not\perp\!\!\!\perp X$, $T \not\perp\!\!\!\perp Y$, $T \perp\!\!\!\perp X|Y$, and $T \perp\!\!\!\perp Y|X$. A related notion is *conditional information equivalence* (Lemeire et al., 2012; Statnikov et al., 2013): X and Y are called to be conditionally information equivalent with respect to T given $Z \subseteq V \setminus X \setminus Y \setminus T$, if $T \not\perp\!\!\!\perp X|Z$, $T \not\perp\!\!\!\perp Y|Z$, $T \perp\!\!\!\perp X|Y \cup Z$, and $T \perp\!\!\!\perp Y|X \cup Z$. Lemeire et al. (2012, pp. 1309–1311) showed that (conditional) information equivalence is one of the two major cases in which *adjacency faithfulness* is violated. Here, the adjacency faithfulness condition (Ramsey et al., 2006; Lemeire et al., 2012) is defined as: if X and Y are adjacent, then $X \not\perp\!\!\!\perp Y|Z$ for any $Z \subseteq V \setminus \{X, Y\}$. Statnikov et al. (2013, p. 503) provided a local version for adjacency faithfulness by focusing on a specific variable.

Here, we employ the information flow metaphor (Cheng et al., 2002) to intuitively explain information equivalence: we can view a BN as a network of information channels, where each node is a *valve* that is either active or inactive; the valves are connected by *information channels*; information can flow through an active valve but not an inactive one; instantiating a node means this valve becomes inactive. We extend this metaphor by viewing a clique of one or more nodes as a valve. In this sense, a CI relation $X \perp\!\!\!\perp Y|Z$ means all the channels between X and Y are cut off by Z and thus the information between X and Y can not flow once Z becomes inactive. When information equivalence occurs, we further extend this information flow metaphor as follows: if X and Y are information equivalent with respect to T given Z , then there exists an *information equivalent valve*, denoted by $\delta_{X,Y;T|Z}$, which connects T and X and connects T and Y ; $\delta_{X,Y;T|Z}$ is active when and only when both X and Y are active. Then, the relation “ X and Y are information equivalent with respect to T given Z ” can be presented in Figure 2.

Information equivalence represents all the possible situations leading to the nonuniqueness of MB. In fact, we can show the following result, which indicates that violating the uniqueness of MB implies the presence of information equivalence. The proof is presented in Section B.

Lemma 1 *The intersection property holds if and only if no information equivalence occurs.* ■

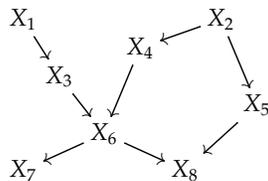


Figure 1: A simple DAG used to illustrate d-separation.

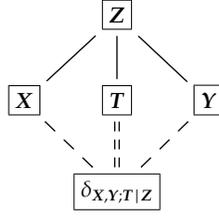


Figure 2: An intuitive illustration for information equivalence.

By analyzing the proof of Lemma 1, we present a relaxed version for the intersection property called *local intersection* as follows.

Definition 2 (Local Intersection) Letting $T \subseteq V$, we say T satisfies the local intersection property, if the following two types of local conditions hold simultaneously: (i) *type-I local condition*: in the case of $|T| \geq 2$, for any disjoint $T_1, T_2 \subseteq T$, there are no disjoint $X, Y \subseteq V \setminus T$ such that T_1 and T_2 are information equivalent with respect to X conditioned on Y ; and (ii) *type-II local condition*: there are no disjoint $X, Y, Z \subseteq V \setminus T$ such that X and Y are information equivalent with respect to T conditioned on Z . ■

Clearly, *intersection* implies *local intersection* but not vice versa, because the former requires no any information equivalence while the latter only requires no information equivalence between the targets and the remaining variables. Here, we give a lemma concerning the uniqueness of MB under the local intersection assumption. The proof is presented in Appendix B.

Lemma 2 For $T \subseteq V$, assume the type-II local condition defined in Definition 2 holds. Then T has a unique MB. ■

To facilitate the identification of information equivalence, Lemeire (2007) ever introduced the notions of target partition (T-partition) and equivalent partition (E-partition), and then provided a relation among information equivalence, T-partition, and E-partition.

- T-partition: The domain, X_{dom} , of X can be partitioned into disjoint subsets $X_{\text{dom}}^{(k)}$ for which $\mathbb{P}(T|\mathbf{x})$ is the same for all $\mathbf{x} \in X_{\text{dom}}^{(k)}$. This is called the T-partition of X_{dom} with respect to T .
- E-partition: A relation $\mathcal{R} \subset X \otimes Y$ defines an E-partition in Y_{dom} to a partition of X_{dom} , if: (i) $\neg(x_2 \mathcal{R} y_1)$ holds for any $x_1, x_2 \in X_{\text{dom}}$ belonging to different partitions and for any $y_1 \in Y_{\text{dom}}$ with $x_1 \mathcal{R} y_1$; and (ii) for every $X_{\text{dom}}^{(k)}$, there exist $x_1 \in X_{\text{dom}}^{(k)}$ and $y_1 \in Y_{\text{dom}}$ such that $x_1 \mathcal{R} y_1$.
- Relation among *information equivalence*, *T-partition*, and *E-partition*: If $T \not\perp X$ and $T \perp Y | X$, then $T \perp X | Y$ (meaning X and Y are information equivalent with respect to T) if and only if the relation $x \mathcal{R} y$ defined by $\mathbb{P}(x, y) > 0$ with $x \in X_{\text{dom}}$ and $y \in Y_{\text{dom}}$ defines an E-partition in Y_{dom} to the T-partition of X_{dom} with respect to T .

The graph shown in Figure 3, originally presented by Statnikov and Aliferis (2010), makes an intuitive illustration on T-partition and E-partition. As seen, $\{1, 2\}$ and $\{3\}$ constitute the T-partition of $A_{\text{dom}} \triangleq \{1, 2, 3\}$ with respect to C ; $\{1, 2\}$ and $\{3\}$ are the E-partition of $B_{\text{dom}} \triangleq \{1, 2, 3\}$ to the T-partition of A_{dom} . Therefore, A and B are information equivalent with respect to C if $C \not\perp A$, since $C \perp B | A$ holds inherently because of the Markov condition.

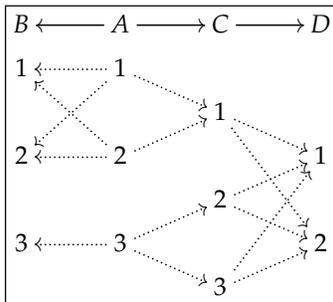


Figure 3: An illustration on T-partition and E-partition: in the DAG “ $B \leftarrow A \rightarrow C \rightarrow D$ ”, all variables take $\{1, 2, 3\}$ except for D taking $\{1, 2\}$, and dotted arrows denote all non-zero conditional probabilities of each variable given its parents.

Finally, the notion of *context-independent information equivalence* given by Statnikov et al. (2013) will be used in Example 2. X and Y are called context-independent information equivalent with respect to T , if X and Y are information equivalent with respect to T given any $Z \subseteq V \setminus X \setminus Y \setminus T$. For this notion, Statnikov et al. (2013) proved the following conclusion: if M is an Mb of T with $X \subseteq M$, and there is some $Y \subseteq V \setminus M \setminus T$ such that X and Y are context-independent information equivalent with respect to T , then $(M \setminus X) \cup Y$ is also an Mb of T .

2.2 Two Typical Algorithms: IAMB and KIAMB

This subsection concisely presents two typical MB discovery algorithms: IAMB (Tsamardinos et al., 2003) and KIAMB (Peña et al., 2007). We select them because of their high adaptability and time efficiency: (i) correctness of IAMB and KIAMB requires only the local composition assumption (Statnikov et al., 2013), while the correctness of the parents and children based algorithms, such as PCMB and the algorithms in the GLL framework, usually requires the faithfulness condition (Peña et al., 2007, Theorem 6; and Aliferis et al., 2010a, Theorem 1); (ii) IAMB and KIAMB are time efficient and thus suitable for the problem of MB for multiple variables, while the parents and children based algorithms have exponential complexities (Aliferis et al., 2010a, pp. 199–200), so they are hard to work when too many variables are involved, such as the problem of MB discovery for multiple variables.

IAMB is an enhanced variant of GS. In 2003, Tsamardinos et al. pointed out that GS uses a static and potentially inefficient heuristic in the growing phase, and then proposed IAMB by employing a dynamic heuristic. Tsamardinos et al. (2003) showed the correctness of IAMB under the faithfulness condition; Peña et al. (2007) relaxed the condition to the composition assumption; Statnikov et al. (2013) further relaxed the condition to the local composition assumption. Algorithm 3 describes the pseudo code for IAMB. See Appendix A for details.

In the algorithm, there is a function f_D (Line 3 of IAMB in Algorithm 3) denoting a heuristic used to measure the association between variables (Tsamardinos et al., 2003; Peña et al., 2007). Two widely used selections for f_D are CMI (Cheng et al., 2002; Tsamardinos et al., 2003) and the negative p -value (Tsamardinos et al., 2006; Aliferis et al., 2010a,b; Statnikov et al., 2013). Also, Yaramakala (2004, p. 41) suggested an equivalent version of the negative p -value. Subsection 4.3 will make a discussion about the ways of practically doing CI tests and the selections for f_D .

KIAMB is a stochastic extension of IAMB. It embeds a randomization parameter $K \in [0, 1]$ which specifies the trade-off between greediness and randomness. If taking $K = 1$, KIAMB reduces to IAMB. Peña et al. (2007) proved the correctness of KIAMB under the composition assumption. By the proof, the local composition assumption is sufficient for KIAMB to be correct. Algorithm 3 describes the pseudo code for KIAMB.

For the case of $|\mathbf{T}| \geq 2$, IAMB and KIAMB can remain correct if strengthening the precondition. We present the correctness of them as follows, without presenting the proof since it is similar to that of the original IAMB and KIAMB (Tsamardinos et al., 2003; Peña et al., 2007; Statnikov et al., 2013). In what follows, we say a CI test for a hypothesis is correct if the statistical decision is correctly made by using a testing method. Subsection 4.3 gives a further discussion on this issue.

Theorem 1 (Correctness of IAMB and KIAMB) *Assume \mathbf{T} satisfies the local composition property, and all CI tests are correct. Then (i) IAMB outputs an MB of \mathbf{T} ; (ii) KIAMB outputs an MB of \mathbf{T} for any $K \in [0, 1]$. ■*

2.3 Motivations

This subsection provides three motivations of this paper.

Let \mathbf{M} be an MB of \mathbf{T} . Then $\mathbb{P}(\mathbf{T} | \mathbf{V} \setminus \{\mathbf{T}\}) = \mathbb{P}(\mathbf{T} | \mathbf{M})$. In other words, all information for predicting \mathbf{T} is carried by \mathbf{M} . Further, \mathbf{M} is a solution to the FS problem, if the algorithm that constructs the prediction model can learn any probability distribution, and the performance metric is strictly decreasing with the mean-squared loss with a preference for smaller subsets (Tsamardinos and Aliferis, 2003, Proposition 3). For this reason, MB for a single variable is sufficient.

However, there are the situations where MB for multiple variables is preferred. This occurs if we need the probability distribution of more than one variables given all the others. Let \mathbf{M}_i be an MB of T_i for $i = 1, 2$. Denoting $\mathbf{T} = \{T_1, T_2\}$, it follows that

$$\mathbb{P}(\mathbf{T} | \mathbf{V} \setminus \mathbf{T}) = \begin{cases} \mathbb{P}(T_1 | \mathbf{M}_1) \mathbb{P}(T_2 | \mathbf{M}_2) & \text{if } T_1 \notin \mathbf{M}_2 \text{ or } T_2 \notin \mathbf{M}_1 \\ \mathbb{P}(T_1, T_2, \mathbf{V} \setminus \mathbf{T}) / \sum_{t_1, t_2} \mathbb{P}(t_1, t_2, \mathbf{V} \setminus \mathbf{T}) & \text{if } T_1 \in \mathbf{M}_2 \text{ and } T_2 \in \mathbf{M}_1 \end{cases}$$

As seen, in the case of $T_1 \in \mathbf{M}_2$ and $T_2 \in \mathbf{M}_1$, the computation is intractable, especially when the dimension is high. Nevertheless, if we have an MB for \mathbf{T} , denoted by \mathbf{M} , then $\mathbb{P}(\mathbf{T} | \mathbf{V} \setminus \mathbf{T}) = \mathbb{P}(\mathbf{T} | \mathbf{M})$ follows immediately, so the problem is simplified greatly. In this sense, it is meaningful to consider the problem of MB for multiple variables.

The second motivation is that we want to know whether the prediction for \mathbf{T} will be affected if the observed values of some variables outside \mathbf{T} (in a new observation) are missing. Denote these missing variables by \mathbf{V}_m . This problem can be considered as follows: find an approximate MB (denoted by \mathbf{M}_m) of \mathbf{T} in $\mathbf{V} \setminus \mathbf{V}_m$ by means of some method, then check if \mathbf{M}_m is an Mb in \mathbf{V} via some criterion (e.g., a criterion based on Lemma 2 given by Statnikov et al., 2013); and finally assert \mathbf{T} will not be affected if the above checking result is “yes”. In this sense, it is also preferred to consider MB for multiple variables.

Figure 4 represents the DAG for the ALARM network (Beinlich et al., 1989), which is well known in the literature. Take $T_1 \triangleq X_{22}$ and $T_2 \triangleq X_{23}$. Then, \mathbf{M}_{T_i} is the unique MB of T_i for $i = 1, 2$ under the faithfulness condition, with

$$\mathbf{M}_{T_1} \triangleq \{X_1, X_4, X_{15}, X_{21}, X_{23}, X_{27}, X_{29}\} \text{ and } \mathbf{M}_{T_2} \triangleq \{X_2, X_{22}, X_{24}, X_{25}, X_{27}, X_{29}\}, \quad (1)$$

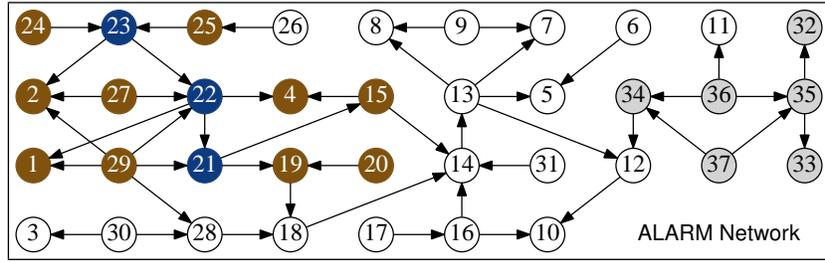


Figure 4: ALARM network (37 nodes and 46 edges): a logical alarm reduction mechanism.

respectively. This leads to intractable computations on the joint probability distribution of T_1 and T_2 given all other variables, consider that $T_1 \in \mathbf{M}_{T_2}$ and $T_2 \in \mathbf{M}_{T_1}$. Further, if the observed values of some variables (e.g., X_j for $j = 32, 33, \dots, 37$) in a new observation are missing, can this observation be used any more for predicting T_1 and T_2 ? Furthermore, we have to face similar problems if three or more target variables are considered.

The third motivation concerns MB discovery algorithms. By Theorem 1, IAMB and KIAMB can be applied to the problem of MB for multiple variables if simply regarding the targets as a multivariate vector, under the strengthened local composition assumption. However, the assumption of local composition imposed on multiple targets may have more occasions to become invalid than imposed on single targets, due to the synergy effect in the sense that neither X nor Y carries information of T but together they contain some information of T (Rauh et al., 2014).

Here is an illustration: considering the BN with the graph in Figure 5 as its DAG, by direct computations using the FullBNT toolbox (Murphy, 2007), we find that

$$\begin{aligned} A \perp\!\!\!\perp C, & & A \perp\!\!\!\perp D, & & \text{and} & & A \perp\!\!\!\perp \{C, D\}; \\ B \perp\!\!\!\perp C, & & B \perp\!\!\!\perp D, & & \text{and} & & B \perp\!\!\!\perp \{C, D\}; \\ \{A, B\} \perp\!\!\!\perp C, & & \{A, B\} \perp\!\!\!\perp D, & & \text{but} & & \{A, B\} \not\perp\!\!\!\perp \{C, D\}. \end{aligned}$$

By this illustration, the idea of applying the existing MB discovery algorithms to multiple targets seems to be practically improper although it is theoretically feasible, because synergy effects may lead to potential inefficiency and even incorrectness. This motivates us to build some algorithms which are resistant to synergy effects and, further, are time efficient.

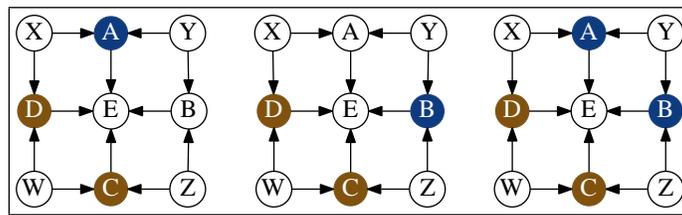


Figure 5: An illustration on synergy effects: each of $\{X, Y, Z, W\}$ takes $\{1, 2\}$ equiprobably; each of $\{A, B, C, D\}$ takes 1 with probabilities p_1, p_2, p_3, p_4 and takes 2 with probabilities $1 - p_1, 1 - p_2, 1 - p_3, 1 - p_4$ given its parents, with $p_4 = p_1 - p_2 + p_3$; E has an arbitrary distribution.

3. Markov Blanket and Markov Boundary for Multiple Variables

This section presents the theoretical results on the problem of Mb and MB for multiple variables when the local intersection property is satisfied and when this property is violated. We study this problem following this way because we are trying to find a suitable approach to transform the problem of Mb and MB from multiple case to single cases, based on which we can build efficient algorithms with high accuracies and low complexities.

3.1 Additivity under Local Intersection

In this subsection, we consider the problem of Mb and MB for multiple variables under the local intersection assumption. We prove Mb and MB possess an ideal property called *additivity*. That is, an Mb of multiple variables can be constructed by simply taking the union of the Mbs of the individual variables and removing the target variables themselves (the same for MB). The results are presented in Theorem 2 and Theorem 3, respectively. Appendix B gives their proofs.

Theorem 2 (Additivity of Mb) *Let (\mathbb{G}, \mathbb{P}) be a BN over V . The following two statements hold:*

- (i) *Let M_i be an Mb of $T_i \subseteq V$ for $i = 1, 2$, and assume $T_1 \cup T_2$ satisfies the local intersection assumption. Then, $(M_1 \cup M_2) \setminus (T_1 \cup T_2)$ is an Mb of $T_1 \cup T_2$.*
- (ii) *Let M_i be an Mb of $T_i \in V$ for $i = 1, \dots, k$, and assume $T \triangleq \{T_1, \dots, T_k\}$ satisfies the local intersection assumption. Then, $\bigcup_{i=1}^k M_i \setminus T$ is an Mb of T . ■*

The additivity property of Mb can be intuitively described by the information flow metaphor (Cheng et al., 2002) using Figure 6: $(M_1 \cup M_2) \setminus (T_1 \cup T_2)$ is enough to cut off all information channels from $T_1 \cup T_2$ to other valves, when no information equivalence associated with $T_1 \cup T_2$ occurs.

Let $T \subseteq V$ be the set of target variables. As we know, in the case of $|T| = 1$ (denoting $T = \{T\}$), the set M_T composed of the parents, children, and spouses of T is an Mb of it (Pearl, 1988), since M_T d-separates T from all other variables. For the case of $|T| \geq 2$ (denoting $T = \{T_1, \dots, T_k\}$), Theorem 2 indicates that the union of all M_{T_i} 's with T_1, \dots, T_k excluded is an Mb of T .

Considering the ALARM network presented in Figure 4, we put $T_1 \triangleq X_{22}$ and $T_2 \triangleq X_{23}$. Then M_{T_i} is an Mb of T_i for $i = 1, 2$, where M_{T_1} and M_{T_2} are defined in (1). Assume $T_{1,2} \triangleq \{T_1, T_2\}$ satisfies the local intersection property. It follows from Theorem 2 that

$$(M_{T_1} \cup M_{T_2}) \setminus T_{1,2} = \{X_1, X_2, X_4, X_{15}, X_{21}, X_{24}, X_{25}, X_{27}, X_{29}\} \triangleq M_{1,2} \quad (2)$$

is an Mb of $T_{1,2}$. Those variables outside $M_{1,2}$ contain no information about $T_{1,2}$ conditioned on $M_{1,2}$ and thereby $\mathbb{P}(T_{1,2} | V \setminus T_{1,2})$ reduces to $\mathbb{P}(T_{1,2} | M_{1,2})$. Further, if the observed values of

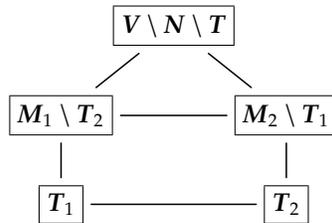


Figure 6: An illustration for additivity of Mb and MB with $T = T_1 \cup T_2$ and $N = (M_1 \cup M_2) \setminus T$.

X_j for $j = 32, \dots, 37$ in a new observation are missing, this observation can still be used without affecting the prediction on $T_{1,2}$. Further, $M_{T_3} \triangleq \{X_{15}, X_{19}, X_{20}, X_{22}, X_{29}\}$ is an Mb of $T_3 \triangleq X_{21}$. Assume $T_{1,2,3} \triangleq \{T_1, T_2, T_3\}$ satisfies the local intersection property. Then Theorem 2 shows

$$M_{T_1} \cup M_{T_2} \cup M_{T_3} \setminus T_{1,2,3} = \{X_1, X_2, X_4, X_{15}, X_{19}, X_{20}, X_{24}, X_{25}, X_{27}, X_{29}\} \triangleq M_{1,2,3} \quad (3)$$

is an Mb of $T_{1,2,3}$.

For additivity of Mb shown in (i) of Theorem 2, we have a useful remark (used to simplify our algorithms in Section 4), based on the fact that if M is an Mb of T then $M \cup M_0$ is also an Mb of T for any $M_0 \subseteq V \setminus M \setminus T$. By the remark, the local intersection assumption for additivity of Mb is not required in some special cases. The proof of this remark is given in Appendix B.

Remark 1 *In the case of either $T_1 \subseteq V \setminus M_2$ or $T_2 \subseteq V \setminus M_1$, the conclusion of (i) in Theorem 2 holds without requiring the local intersection assumption.* ■

Theorem 2 shows the additivity of Mb. A natural idea is to wonder if additivity is possessed by MB. Theorem 3 affirms this. Appendix B provides the proof. Note that the statements about the uniqueness of MB in this theorem follow from Lemma 2.

Theorem 3 (Additivity of MB) *Let (\mathbb{G}, \mathbb{P}) be a BN over V . The following two statements hold:*

- (i) *Assume $T_1 \cup T_2$ satisfies the local intersection assumption. Let M_i be the unique MB of T_i for $i = 1, 2$. Then, $(M_1 \cup M_2) \setminus (T_1 \cup T_2)$ is the unique MB of $T_1 \cup T_2$.*
- (ii) *Assume $T \triangleq \{T_1, \dots, T_k\}$ satisfies the local intersection assumption. Let M_i be the unique MB of T_i for $i = 1, \dots, k$. Then, $\bigcup_{i=1}^k M_i \setminus T$ is the unique MB of T .* ■

According to Theorem 3, $M_{1,2}$ defined in (2) is not only an Mb but also the unique MB of $T_{1,2}$ in the ALARM network if the faithfulness condition is satisfied. Further, $M_{1,2,3}$ defined in (3) is the unique MB of $T_{1,2,3}$.

3.2 Theoretical Results in the General Case

Let (\mathbb{G}, \mathbb{P}) be a BN over V , and assume $T_i \subseteq V$ with $|T_i| \geq 1$ has an Mb or MB, M_i , for $i = 1, 2$. Denote $T = T_1 \cup T_2$ and $N = (M_1 \cup M_2) \setminus T$. In the case that M_i is an MB of T_i , Theorem 3 reveals that N is an MB of T if T satisfies the local intersection assumption. However, when the local intersection assumption does not hold (meaning information equivalence occurs, as Lemma 1 shows), N may be no longer an MB of T , due to one of the following reasons: (i) N may be an Mb but it may not possess minimality, as shown by Example 2; (ii) N may be insufficient to shield T_1 and T_2 from all other variables, so it is no longer an Mb in this case, and some extra variables are required to enter into N . Example 1 provides an illustration.

For the first case, we need only to optimize N by simply removing redundant variables from N ; however, for the second case, the additivity property of MB is thoroughly broken, and the problem of constructing an MB for T based on M_1 and M_2 becomes complex. On the one hand, there are some variables in $V \setminus N \setminus T$ needing to enter into N ; on the other hand, there may be some variables in N becoming redundant once some new members supplement N . What we concern are which variables should enter into N and how we find them.

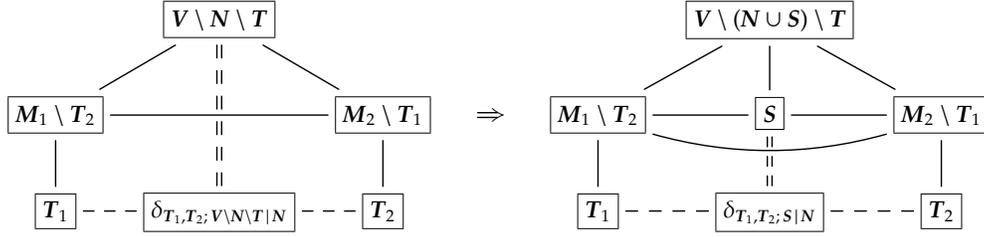


Figure 7: An illustration for the case of violating additivity of Mb and MB, caused by information equivalence.

Assume N is no longer an Mb of T . Then, it is easily shown that

$$\begin{aligned} V \setminus N \setminus T \not\perp\!\!\!\perp T_1 | N, \quad \text{but} \quad V \setminus N \setminus T \perp\!\!\!\perp T_1 | N \cup T_2, \\ V \setminus N \setminus T \not\perp\!\!\!\perp T_2 | N, \quad \text{but} \quad V \setminus N \setminus T \perp\!\!\!\perp T_2 | N \cup T_1. \end{aligned}$$

That is, T_1 and T_2 contain equivalent information about $V \setminus N \setminus T$ given N . See Figure 7 for an illustration: the valves $M_1 \setminus T_2$ and $M_2 \setminus T_1$ can not cut off all information channels between T and $V \setminus N \setminus T$, because some information can flow through $\delta_{T_1, T_2; V \setminus N \setminus T | N}$, an information equivalent valve of T_1 and T_2 with respect to $V \setminus N \setminus T$ given N . In other words, T_1 and T_2 may exchange information directly; besides, they also share the equivalent information about $V \setminus N \setminus T$. This indicates we should continue to turn off some valves, $S \subseteq V \setminus N \setminus T$, besides $M_1 \setminus T_2$ and $M_2 \setminus T_1$ such that T_1 and T_2 no longer exchange information through external valves and thus such that T has no information exchange with remaining valves.

This analysis motivates us to give the following definition:

Definition 3 *With the notations above, we call $S (\subseteq V \setminus N \setminus T)$ a Markov blanket supplementary (MbS) (of T to N), if $N \cup S$ is an Mb of T . Further, a Markov boundary supplementary (MBS) is any MbS such that none of its proper subsets is an MbS. ■*

In what follows, we give the properties of MbS and MBS.

Theorem 4 *Assume $S \subseteq V \setminus N \setminus T$. Then, the following statements are equivalent:*

- (i) S is an MbS;
- (ii) $\mathbb{I}(T_1; T_2 | N \cup S) = \min_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T_1; T_2 | N \cup S')$;
- (iii) $\mathbb{I}(T; S | N) = \max_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T; S' | N)$;
- (iv) $N \cup S$ is an Mb of T_1 in $V \setminus T_2$ (or $N \cup S$ is an Mb of T_2 in $V \setminus T_1$).

In addition, if S is an MbS, then it is also an MBS if and only if $T_1 \not\perp\!\!\!\perp Y | N \cup (S \setminus \{Y\})$ or $T_2 \not\perp\!\!\!\perp Y | N \cup (S \setminus \{Y\})$ holds for any $Y \in S$. ■

The proof of this theorem is presented in Appendix B.

As seen, (ii) and (iii) of Theorem 4 explain the implication of MbS that the information flow metaphor illustrates in Figure 7: finding an MbS is equivalent to turning off some valves such that T_1 and T_2 no longer exchange information through external valves, or equivalent to finding all remaining equivalent information contained by T_1 and T_2 ; (iv) and the property of MBS provide a practical way of building MBS discovery algorithms.

Here, we use an example to demonstrate the notions of MbS and MBS and their properties.

Example 1 Consider the BN (\mathbb{G}, \mathbb{P}) over $V = \{A, B, C, D\}$ presented in Figure 8, in which $A, B,$ and C take $\{1, 2, 3\}$ while D takes $\{1, 2\}$. Put $T = \{T_1, T_2\}$, $N = (M_1 \cup M_2) \setminus T = \emptyset$, and $S = \{C\}$, $S_0 = \{C, D\}$ with $T_1 = A$, $T_2 = B$, $M_1 = \{B\}$, $M_2 = \{A\}$. Using the theory of information equivalence (Lemeire, 2007), we can show the following results (see Appendix B for the proofs):

- (i) M_1 is an MB of T_1 in V : $\mathbb{I}(A; C, D | B) = 0$ and $\mathbb{I}(A; C, D) > 0$;
- (ii) M_2 is an MB of T_2 in V : $\mathbb{I}(B; C, D | A) = 0$ and $\mathbb{I}(B; C, D) > 0$;
- (iii) $N \cup S$ is an Mb of T in V , so S is an Mbs: $\mathbb{I}(A, B; D | C) = 0$;
- (iv) $\mathbb{I}(T_1; T_2 | N \cup S) = \min_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T_1; T_2 | N \cup S')$, because of $\mathbb{I}(A; B | C) = \mathbb{I}(A; B | C, D)$, $\mathbb{I}(A; B | C) \leq \mathbb{I}(A; B | D)$, and $\mathbb{I}(A; B | C) \leq \mathbb{I}(A; B)$;
- (v) $\mathbb{I}(T; S | N) = \max_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T; S' | N)$;
- (vi) $N \cup S$ is an MB of T_1 in $V \setminus \{T_2\}$: $\mathbb{I}(A; C, D) > 0$ and $\mathbb{I}(A; D | C) = 0$;
- (vii) $N \cup S$ is an MB of T_2 in $V \setminus \{T_1\}$: $\mathbb{I}(B; C, D) > 0$ and $\mathbb{I}(B; D | C) = 0$;
- (viii) S is an MBS; S_0 is an Mbs (not an MBS): $\mathbb{I}(A, B; C, D) > 0$ and $\mathbb{I}(A; B | C, D) = \mathbb{I}(A; B | C)$. ■

By Example 1, A and B share the equivalent information about C , so turning off the valve A (or B) means cutting off all the channels from B (or A) to C . This is why they can screen off each other from C . However, A and B lose the shield if they are integrated into a whole. In this case, we have to turn C off such that A and B no longer exchange information through external valves. This example reveals that an MBS is a minimal set of variables, $S \subseteq V \setminus N \setminus T$, such that T_1 and T_2 contain no equivalent information about the remaining variables given $N \cup S$.

When finding an MBS, S , and letting the variables in S supplement N , there may be some variables in N becoming redundant. In addition, N may be redundant even before supplementing S . Example 2 gives an illustration. For both cases, we need to remove the redundant variables.

Example 2 Consider the BN presented in Figure 9, in which any one variable from $\{A, B, C\}$ and another from $\{D, E, F\}$ (denoted by X and Y , respectively) contain context-independent equivalent information about G (see Statnikov et al., 2013, Example 3). Then, $\{X, Y\}$ is an MB of G . Put now $T_1 = \{C, F\}$ and $T_2 = \{G\}$, and take $M_1 = \{B, E\}$ and $M_2 = \{B, D\}$. Note that $T_1 \subseteq V \setminus M_2$ (and also $T_2 \subseteq V \setminus M_1$). It concludes that $N = \{B, D, E\}$ is not an MB but only an Mb of $T_1 \cup T_2$, since its proper subset $\{B, E\}$ is also an Mb (and also an MB) of $T_1 \cup T_2$. This shows why the process of refining N is necessary. ■

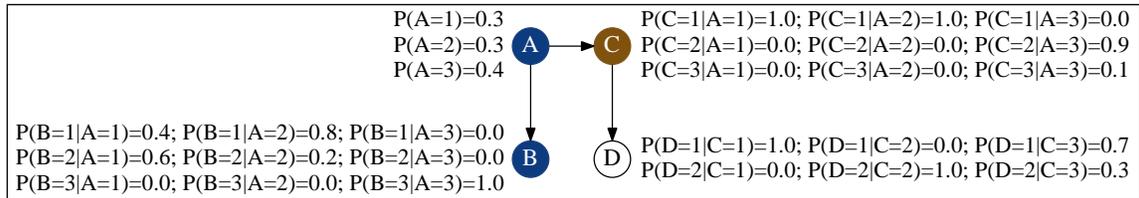


Figure 8: BN (\mathbb{G}, \mathbb{P}) : \mathbb{P} is a joint probability distribution over $V = \{A, B, C, D\}$ with each variable taking values $\{1, 2, 3\}$ except for D taking $\{1, 2\}$; \mathbb{G} is a DAG over V .

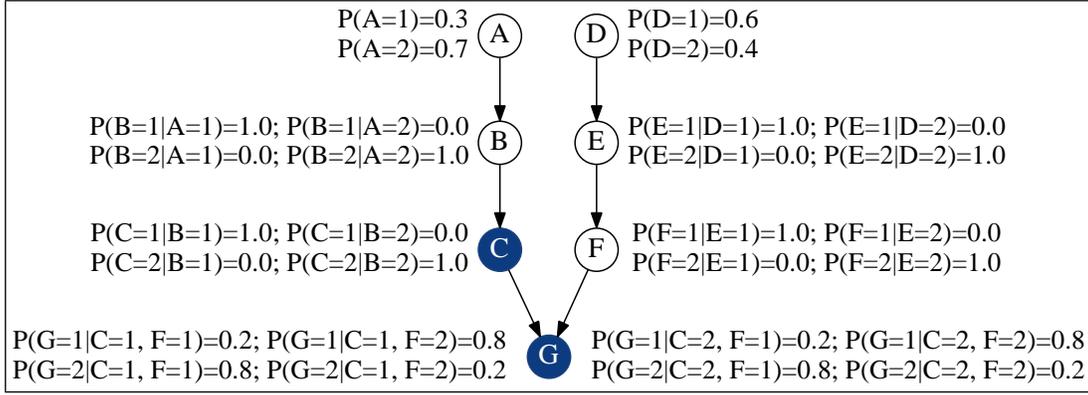


Figure 9: BN (\mathbb{G}, \mathbb{P}) : \mathbb{P} is a joint probability distribution over $V = \{A, B, C, D, E, F, G\}$ with all variables taking values $\{1, 2\}$, \mathbb{G} is a DAG with the variables in V as its nodes.

3.3 An Alternative Approach

Before building MB discovery algorithms for multiple targets, this subsection concisely presents an alternative approach to the *additivity based* and *MBS based* methods. In Section 5, we will apply this method as an FS strategy in multi-class prediction problems.

Let $T \triangleq \{T_1, \dots, T_k\}$ be the targets of interest and T be T 's merged version, taking values $\{1, \dots, t\}$ with $t \geq 3$. This procedure transforms the MB discovery for multiple targets T into the MB discovery for single target T , so all the existing MB discovery algorithms can be employed theoretically if the required conditions are satisfied. However, if t is large, selecting features of T or T directly will be difficult. Subsection 3.1 and Subsection 3.2 provide a way of solving this problem in different situations. In this case, an alternative strategy is to further convert T into a set of dummy variables denoted by $\{T_j^{(d)}\}_{j=1}^t$, where $T_j^{(d)}$ is a 0-1 variable defined as

$$T_j^{(d)} = \begin{cases} 1, & \text{if } T = j \\ 0, & \text{if } T \neq j \end{cases}$$

This transformation produces a multiple-target $T^{(d)} \triangleq (T_1^{(d)}, \dots, T_t^{(d)})$. Clearly, T , T , and $T^{(d)}$ have the same MBs. In what follows, we show the MB of $T^{(d)}$ can be derived by simply taking the union of MBs of $T_1^{(d)}, \dots, T_t^{(d)}$ and then removing the redundant variables in an efficient way. The proof will be given in Appendix B.

Theorem 5 Let M_j be an MB of $T_j^{(d)}$ in $V \setminus T$ for $j = 1, \dots, t$. Then, $M \triangleq \cup_{j=1}^t M_j$ is an Mb of T . Further, M is an MB of T iff for any $X \in M$ there is some j such that $T_j^{(d)} \not\perp X \mid M \setminus \{X\}$. ■

For why this transformative method is efficient, the fourth concluding remark in Section 7 will make a brief explanation.

4. Algorithms

This section builds MB discovery algorithms for multiple targets, $\{T_1, \dots, T_k\} \triangleq T$.

Let \mathbb{A} be an MB discovery algorithm, assumed to perform well when used to discover an MB for a single target. In this paper, we employ IAMB and KIAMB as \mathbb{A} . Clearly, \mathbb{A} can be directly used to find an MB for T if simply regarding T as the input of \mathbb{A} . Usually, this will lead to low accuracies and high complexities.

By Theorem 4, the MB discovery problem for multiple targets can be translated equivalently into a number of MB discovery problems for single targets, according to the following way: (i) use \mathbb{A} to find an MB of T_i in V for $i = 1, \dots, k$, denoted by M_i ; (ii) find an Mb of T_2 in $V \setminus \{T_1\}$ based on $(M_1 \cup M_2) \setminus \{T_1, T_2\}$, and then get an MB of $\{T_1, T_2\}$, written as $M_{1,2}$; (iii) find an Mb of T_3 in $V \setminus \{T_1, T_2\}$ based on $(M_{1,2} \cup M_3) \setminus \{T_1, T_2, T_3\}$, and then get an MB of $\{T_1, T_2, T_3\}$, written as $M_{1,2,3}$; (iv) the rest can be done in a similar manner. Following this way, the input of \mathbb{A} for each use is a single variable, so this idea successfully avoids assigning an multivariate input to \mathbb{A} . Note that in the above process the equivalent information is extracted in a stepwise manner.

4.1 IAMBS and KIAMBS

Let (\mathbb{G}, \mathbb{P}) be a BN over V , and assume $T_i \subseteq V$ with $|T_i| \geq 1$ has an MB, M_i , for $i = 1, 2$. Denote $N = (M_1 \cup M_2) \setminus (T_1 \cup T_2)$. This subsection presents the algorithms for discovering an MB of $T_1 \cup T_2$. To design one such algorithm, we note that there may be some variables in N becoming redundant once an MBS, S , is supplemented. Therefore, we need to first find S by setting N as a whitelist in \mathbb{A} and then refine N .

Applying this idea to IAMB and KIAMB, we obtain two algorithms called IAMBS and KIAMBS, in which “S” refers to as “supplementary”. Their pseudo codes are presented in Algorithm 1. In order to differentiate these two algorithms, we set K in the KIAMBS algorithm as $K \in [0, 1)$. It is mentioned here that S_2 is a random subset of S_1 with size $\max\{1, \lfloor |S_1| \cdot K \rfloor\}$ in Line 5 of KIAMBS. As seen, these two algorithms first find an MbS, S , in the growing phase and then refine S and N in sequence in the shrinking phase.

For example, based on a data set drawn from the BN in Example 1, the unique MB, $\{C\}$, of $\{A, B\}$ can be discovered by calling IAMBS or KIAMBS only once.

Theorem 1 presents the correctness of IAMB and KIAMB under the assumption that $T_1 \cup T_2$ satisfies the local composition property. The theorem below shows IAMBS and KIAMBS are correct if T_2 (instead of $T_1 \cup T_2$) satisfies the local composition property. Appendix B gives the proof.

Theorem 6 (Correctness of IAMBS and KIAMBS) *Assume that T_2 satisfies the local composition property, and that all CI tests are correct. Then (i) IAMBS outputs an MB of $T_1 \cup T_2$; (ii) KIAMBS outputs an MB of $T_1 \cup T_2$ for any $K \in [0, 1)$. ■*

The following remark presents a relation among local intersection, local composition, and the adjacency faithfulness condition, under the *orientation faithfulness condition*. The proof is given in Appendix B. Here, the orientation faithfulness condition (Ramsey et al., 2006; Lemeire et al., 2012) is defined as: for any $X, Y, Z \in V$ such that X and Z are adjacent to Y but X is not adjacent to Z , (i) if $X \rightarrow Y \leftarrow Z$, then $X \not\perp Z \mid W$ holds for any $W \subseteq V \setminus \{X, Z\}$ with $Y \in W$; (ii) otherwise, $X \perp Z \mid W$ holds for any $W \subseteq V \setminus \{X, Y, Z\}$.

Remark 2 *The following two statements hold: (a) violating local intersection implies violating adjacency faithfulness; (b) under the orientation faithfulness condition, violating local composition at the end of the first phase of IAMB or KIAMB or IAMBS or KIAMBS means violating adjacency faithfulness. ■*

In addition, the lemma below is useful to explain the succeeding remarks.

Lemma 3 (a) *If there is $P \subseteq M_1 \setminus T_2$ such that $T_1 \perp\!\!\!\perp P | (N \setminus P) \cup T_2$, then $(N \setminus P) \cup T_2$ is an Mb of T_1 ; (b) *If there is $Q \subseteq N \setminus P$ such that $T_1 \perp\!\!\!\perp Q | (N \setminus P \setminus Q) \cup T_2$ and $T_2 \perp\!\!\!\perp Q | (N \setminus P \setminus Q) \cup T_1$, then $(N \setminus P \setminus Q) \cup T_2$ is an Mb of T_1 , and $(N \setminus P \setminus Q) \cup T_1$ is an Mb of T_2 . ■**

For Algorithm 1, we have three remarks below:

- (i) In these two algorithms, the two CI tests for adding members to S and for refining S are based on T_2 instead of T , while the CI test for refining N is based on T instead of T_2 . (a) For the first two CI tests, T_2 can be replaced with T without affecting the correctness of the algorithms, since $T_2 \perp\!\!\!\perp X | N \cup S' \Leftrightarrow T \perp\!\!\!\perp X | N \cup S'$ holds for any $S' \subseteq V \setminus N \setminus T$ and $X \subseteq V \setminus (N \cup S') \setminus T$. However, if we replace T_2 with T , the resulting algorithms will need much longer time to run. This is why we use T_2 in stead of T in these two places. (b) For the third CI test, T can not be replaced with T_2 , because of $T_2 \perp\!\!\!\perp X | (N \setminus X) \cup S \not\Rightarrow T \perp\!\!\!\perp X | (N \setminus X) \cup S$.
- (ii) According to Remark 2, there may be some situations in which both local intersection and local composition are simultaneously violated. In this case, IAMBS and KIAMBS may not

Algorithm 1: IAMBS and KIAMBS	
<p>Procedure: $M \leftarrow \text{IAMBS}(D; T_1, T_2; M_1, M_2)$ Input: a data matrix D; two sets of targets T_1 and T_2; an MB M_i of T_i for $i = 1, 2$. Output: an MB, M, of $T \triangleq T_1 \cup T_2$.</p> <p>//Forward: Growing Phase</p> <pre> 1 $S \leftarrow \emptyset$ 2 while S has changed do 3 $M \leftarrow N \cup S$ 4 $Y \leftarrow \arg \max_{X \in V \setminus M \setminus T} f_D(T_2; X M)$ 5 if $T_2 \not\perp\!\!\!\perp Y M$ then 6 $S \leftarrow S \cup \{Y\}$ 7 end 8 end //Backward: Shrinking Phase 9 foreach $X \in S$ do 10 if $T_2 \perp\!\!\!\perp Y N \cup (S \setminus \{Y\})$ then 11 $S \leftarrow S \setminus \{Y\}$ 12 end 13 end 14 foreach $Y \in N$ do 15 if $T \perp\!\!\!\perp Y (N \setminus \{Y\}) \cup S$ then 16 $N \leftarrow N \setminus \{Y\}$ 17 end 18 end 19 return $M \leftarrow N \cup S$ </pre>	<p>Procedure: $M \leftarrow \text{KIAMBS}(D; T_1, T_2; M_1, M_2; K)$ Input: Besides $\{D, T_i, M_i\}$, $K \in [0, 1)$ is a randomization parameter. Output: an MB, M, of $T \triangleq T_1 \cup T_2$.</p> <p>//Forward: Growing Phase</p> <pre> 1 $S \leftarrow \emptyset$ 2 while S has changed do 3 $M \leftarrow N \cup S$ 4 if $S_1 \leftarrow \{X \in V \setminus M \setminus T : T_2 \not\perp\!\!\!\perp X M\} \neq \emptyset$ then 5 $Y \leftarrow \arg \max_{X \in S_2} f_D(T_2; X M)$ 6 $S \leftarrow S \cup \{Y\}$ 7 end 8 end //Backward: Shrinking Phase 9 foreach $X \in S$ do 10 if $T_2 \perp\!\!\!\perp Y N \cup (S \setminus \{Y\})$ then 11 $S \leftarrow S \setminus \{Y\}$ 12 end 13 end 14 foreach $Y \in N$ do 15 if $T \perp\!\!\!\perp Y (N \setminus \{Y\}) \cup S$ then 16 $N \leftarrow N \setminus \{Y\}$ 17 end 18 end 19 return $M \leftarrow N \cup S$ </pre>

correctly work. Specifically, the *violation of local intersection* means T_1 and T_2 contain equivalent information about $V \setminus N \setminus T$ given N ; while the *violation of local composition* indicates not all equivalent information are successfully extracted by N . Let P and Q be defined as in Lemma 3, and assume $P \cup Q \neq \emptyset$. Then, it can be shown that T_1 and T_2 contain equivalent information about $P \cup Q$ given $N \setminus (P \cup Q)$. This means some equivalent information about $P \cup Q$ shared by T_1 and T_2 conditioned on $N \setminus (P \cup Q)$ may mask some equivalent information about $V \setminus N \setminus T$ contained by T_1 and T_2 conditioned on N . This may be why not all equivalent information can be extracted by N . According to this analysis, a potential remedy is to run IAMBS or KIAMBS by replacing N with a superset of $N \setminus (P \cup Q)$ that is a subset of N .

- (iii) By Remark 1, if $T_1 \subseteq V \setminus M_2$ or $T_2 \subseteq V \setminus M_1$, N must be an Mb of T , so Lines 2~14 of IAMBS and KIAMBS can be omitted. In this case, however, it is still necessary to refine N , because N may not possess minimality. Example 2 illustrates this necessity.

In addition, another problem that we concern is whether we can refine N before seeking S and, if this is the case, which variables in N can be removed directly. We consider this problem because any redundant variable in N can lead to unnecessary inaccuracies when using N as a part of the conditional set in practical computations. Lemma 3 indicates we can do like this. However, to avoid the danger of missing the information about $P \cup Q$ (this occurs if the equivalent information involved in $P \cup Q$ given $N \setminus P \setminus Q$ is different in some sense from any part of the equivalent information involved in $V \setminus N \setminus T$ given N), we recommend to first search the members of S in $V \setminus N \setminus T$ and then check if some variables in $P \cup Q$ are necessary to enter into S when implementing Lines 2~8 of IAMBS and KIAMBS. Note that this will increase the total running time.

4.2 MIAMB and MKIAMB

In this subsection, we present two multivariate Markov boundary discovery algorithms, called MIAMB and MKIAMB, respectively.

Let $\{T_1, \dots, T_k\} \in V$ with M_i as its an MB for $i = 1, \dots, k$. If the local intersection property is satisfied, Theorem 3 shows $\bigcup_{i=1}^k M_i \setminus T$ is an MB of $T \triangleq \{T_1, \dots, T_k\}$. Otherwise, M may be no longer an MB. In this case, we use MIAMB or MKIAMB to seek an MB for T . Given an ordering of T_1, \dots, T_k , saying $\tau \triangleq \{i_1, \dots, i_k\}$, which determines the priorities of the variables in T entering into the queue whose an MB will be sought in the current step, we denote an MB of $\{T_{i_1}, \dots, T_{i_\ell}\} \triangleq T_{i_\ell}^*$ by $M_{i_j}^*$.

With these notations, MIAMB and MKIAMB are pseudo-coded in Algorithm 2. Their correctness, shown by Theorem 7, is a direct consequence of Theorem 1 and Theorem 6. As seen, MIAMB or MKIAMB uses the following stepwise idea: it first finds an MB of two targets $\{T_{i_1}, T_{i_2}\} = \{T_{i_1}\} \cup \{T_{i_2}\}$, and then finds an MB of three targets $\{T_{i_1}, T_{i_2}, T_{i_3}\} = \{T_{i_1}, T_{i_2}\} \cup \{T_{i_3}\}$; the rest can be done in a similar manner until all the k target variables are considered.

Theorem 7 (Correctness of MIAMB and MKIAMB) *Assume that T_i satisfies the local composition property for $i = 1, \dots, k$, and that all CI tests are correct. Denote $T \triangleq \{T_1, \dots, T_k\}$. Then (i) MIAMB outputs an MB of T ; (ii) MKIAMB outputs an MB of T for any $K \in [0, 1]$. ■*

As we know, for any real data, those preconditions (such as faithfulness or local composition) required by a learning algorithm are hard to hold exactly. However, our algorithms can be seen as an improvement over earlier methods. Specifically, IAMB/KIAMBS algorithms require faithfulness or

local composition for multiple targets when used for MB discovery of multiple targets, while our MIAMB/MKIAMB only need local composition for single targets, which may be more close to real situations than faithfulness or local composition for multiple targets.

For MIAMB or MKIAMB, an ordering τ is set in Algorithm 2 mainly because different orderings may lead to different computational complexities. In Subsection 4.4, we will make a complexity analysis about the algorithms, based on which we present a feasible way of selecting τ , under the expectation that our algorithms should be run as quickly as possible. When $|T| = 2$, however, τ is not necessary.

Besides MIAMB/MKIAMB algorithms (which are *MBS based*), we can consider additivity based (Theorem 3) and dummy variables based (Theorem 5) algorithms: (a) the *additivity based* MIAMB or MKIAMB simply takes the union of outputs of IAMB/KIAMB with respect to all single targets as the output; its correctness requires the conditions in Theorem 7 plus Theorem 3; and (b) the *dummy (variables based)* MIAMB/MKIAMB takes the union of the outputs of IAMB/KIAMB with respect to every dummy variable and removes redundant variables; its correctness requires the same condition as in Theorem 7. Throughout this paper, unless specified, MIAMB/MKIAMB denote the MBS based algorithms.

4.3 A Discussion on CI Test

As argued by Aliferis et al. (2010a, p. 200), the quality of an MB discovery algorithm highly depends on the selected CI testing methods. In this subsection, we discuss the ways of practically doing CI tests. Usually, the Pearson's X^2 test or the log-likelihood ratio G^2 test can be employed for this purpose (Yaramakala, 2004; Bromberg and Margaritis, 2009; Aliferis et al., 2010b; Statnikov et al., 2013). Here, the X^2 statistic and the G^2 statistic have the same asymptotic χ^2 distribution. We can also use some experimental testing methods such as the Akaike information criterion-based test (Cressie and Read, 1989; Scutari, 2010).

Algorithm 2: MIAMB and MKIAMB	
<p>Procedure: $M \leftarrow \text{MIAMB}(D; T; \tau)$</p> <p>Input: a data matrix D; a target set $T \triangleq \{T_1, \dots, T_k\}$; and an ordering $\tau \triangleq \{i_1 \dots, i_k\}$.</p> <p>Output: an MB, M, of T.</p> <p>// MIAMB: $M \leftarrow \text{MIAMB}(D; T; \tau)$</p> <p>1 for $\ell \leftarrow 1$ to k do</p> <p>2 $M_{i_\ell} \leftarrow \text{IAMB}(D; \{T_{i_\ell}\})$</p> <p>3 end</p> <p>4 for $\ell \leftarrow 2$ to k do</p> <p>5 $M_{i_\ell}^* \leftarrow \text{IAMBS}(D; T_{i_{\ell-1}}^*, \{T_{i_\ell}\}; M_{i_{\ell-1}}^*, M_{i_\ell})$</p> <p>6 end</p> <p>7 return $M \leftarrow M_{i_k}^*$</p>	<p>Procedure: $M \leftarrow \text{MKIAMB}(D; T; K; \tau)$</p> <p>Input: a data matrix D; a target set T; a randomization parameter $K \in [0, 1)$; and an ordering τ.</p> <p>Output: an MB, M, of T.</p> <p>// MKIAMB: $M \leftarrow \text{MKIAMB}(D; T; K; \tau)$</p> <p>1 for $\ell \leftarrow 1$ to k do</p> <p>2 $M_{i_\ell} \leftarrow \text{KIAMB}(D; \{T_{i_\ell}\}; K)$</p> <p>3 end</p> <p>4 for $\ell \leftarrow 2$ to k do</p> <p>5 $M_{i_\ell}^* \leftarrow$ $\text{KIAMBS}(D; T_{i_{\ell-1}}^*, \{T_{i_\ell}\}; M_{i_{\ell-1}}^*, M_{i_\ell}; K)$</p> <p>6 end</p> <p>7 return $M \leftarrow M_{i_k}^*$</p>

Recall that we are dealing with the MB discovery problem for *multiple* target variables. When the target set, namely T , contains only a few variables (e.g., 1 or 2), the X^2 test or the G^2 test performs quite well in most situations. Unfortunately, when T contains too many variables (e.g., 5 or 6 or even more), X^2 or G^2 may not work well due to the overmany degrees of freedom. See Appendix C for a detailed discussion. In fact, as Cochran (1954, p. 420) recommended about the working rules for X^2 (also applicable to G^2), these two testing methods are unreliable if more than 20% of the cells in contingency tables have an expected count of less than 5 data points; however, such cases frequently arise in practice (Bromberg and Margaritis, 2009; Yaramakala, 2004).

Many authors have considered improving X^2 and G^2 by adjusting the statistics. Lawley (1956) showed that such tests can be improved by multiplying with a suitable scale factor; Hosmane (1986, 1987, 1990) and the pioneer scholars recommended the following two adjustment procedures (i) replace zero observed counts by a positive constant, leaving nonzero counts intact; and (ii) add a positive constant to all the observed counts. Brin et al. (1997) and Silverstein et al. (1998) used two heuristic “solutions” to the problem of low expected counts as follows: (i) simply ignore these cells when calculating X^2 or G^2 ; and (ii) use what is called *contingency table support* (CT-support): a set of items S has CT-support s at the $t\%$ level if at least $t\%$ of the cells in the contingency table for S have value s . Aliferis et al. (2010b) considered a similar heuristic called *heuristic power size*, which denotes the smallest sample size per cell in the contingency table of a reliable CI test.

The above ideas can lead to improvements on X^2 and G^2 to varying degrees if the dimensions are not very high. However, when working on the MB discovery problem for multiple targets, we need more suitable methods to do CI tests. For this reason, we suggest the following *practical operation*: when $|T| \leq 2$, we can (i) use X^2 or G^2 or their variants mentioned above to do CI tests; otherwise, we consider the following testing method: (ii) use CMI and an experimental threshold, ε , to make statistical decisions as Cheng et al. (2002) did, in the sense that $\mathbb{I}_D(X; Y | Z) \geq \varepsilon$ asserts $X \not\perp Y | Z$ while $\mathbb{I}_D(X; Y | Z) < \varepsilon$ concludes $X \perp Y | Z$, where $\varepsilon \triangleq |T|^{a_1} \cdot \frac{100a_2}{n} \cdot \log_2 v$ is related to the sample size, the average number of values that each variable takes, and the number of targets (denoted by n , v , and $|T|$, respectively), in which a_1 and a_2 are two adjusting factors ($a_1 = 0.5$ and $a_2 \in (0.1, 0.5)$ are recommended). The association function, f_D , can be selected as

$$f_D(X; Y | Z) = \mathbb{I}_D(X; Y | Z) \triangleq f_D^{(2)}(X; Y | Z). \quad (4)$$

Besides this experimental method, we can (iii) improve X^2 or G^2 by adjusting the number of the theoretical degrees of freedom.

For the above (iii), to be clear, we consider the G^2 statistic, $G^2(X; Y | Z) \triangleq 2n \cdot \mathbb{I}_D(X; Y | Z)$, which approximates to the chi-square variate with $r \triangleq (r_X - 1)(r_Y - 1)r_Z$ degrees of freedom, namely $\chi^2(r)$, where r_ξ represents the number of configurations for ξ (de Campos, 2006, p. 2158). Denote the p -value by

$$p(X; Y | Z) = \mathbb{P}\{\chi^2(r) \geq G^2(X; Y | Z)\}.$$

Then, the G^2 test asserts $X \perp Y | Z$ if $p(X; Y | Z) > \alpha$ for a significance level α , and concludes $X \not\perp Y | Z$ if $p(X; Y | Z) \leq \alpha$. In this paper, α is set to be 0.05. Aliferis et al. (2010a, pp. 200–201) provided a further discussion about this. Accordingly, the *negative p -value* is used as the association function, f_D , as Tsamardinos et al. (2006), Aliferis et al. (2010a,b), and Statnikov et al. (2013) did:

$$f_D(X; Y | Z) = -p(X; Y | Z) = -\mathbb{P}\{\chi^2(r) \geq G^2(X; Y | Z)\} \triangleq f_D^{(1)}(X; Y | Z). \quad (5)$$

Replace the theoretical value of r in $p(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$ with its a damped version of the form

$$g_{n,\kappa}(r) \triangleq r \left(1 - e^{-\frac{n}{\kappa r}} \right), \quad (6)$$

where $\kappa > 0$ is a constant, based on which $\frac{n}{\kappa}$ measures the amount of valid cells that n sample instances can support. For convenience, we will call the resulted p -value, denoted by $p_g(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$ instead of $p(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$, and the resulted testing method to be the *damped p -value* and the *damped log-likelihood ratio test* (or damped G^2 test). Further, we use the the following association function:

$$f_D(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = -p_g(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = -\mathbb{P}\{\chi^2(g_{n,\kappa}(r)) \geq G^2(\mathbf{X}; \mathbf{Y} | \mathbf{Z})\} \triangleq f_D^{(3)}(\mathbf{X}; \mathbf{Y} | \mathbf{Z}). \quad (7)$$

In Appendix C, we will provide the details for this damping procedure, and give some numerical illustrations about its reasonability. Clearly, the damped G^2 test approximately degenerates into the ordinary G^2 test when taking κ as a very small positive number.

4.4 Complexity Analysis

In the following, we analyze the computational complexities of the four algorithms: IAMB, KIAMB, MIAMB, and MKIAMB. Usually, the number of CI tests can be employed to measure the complexity of a CI-based MB discovery algorithm (Tsamardinos et al., 2003, 2006; Aliferis et al., 2010a), considering there exists efficient implementations of the CMI-based test or the association computation taking time $O(n \log n)$ if the conditional set is small. However, Aliferis et al. (2010a) also mentioned that the running time, denoted by $t_{n,q}$, for computing per CMI-based statistic is linear to the sample size, n , and exponential to the number, q , of variables in the conditional set. This means we should take $t_{n,q}$ into account, not simply using $O(n \log n)$ to measure the complexity.

Assume we are seeking an MB for $\mathbf{T} \triangleq \{T_1, \dots, T_k\}$ according to the ordering τ . Without loss of generality, we assume $\tau = \{1, \dots, k\}$. Consider the case of $k = 2$. Suppose \mathbf{M}_i is an MB of T_i with $|\mathbf{M}_i| = m_i \geq 1$, and \mathbf{S} is an MBS for $\mathbf{N} \triangleq \mathbf{M}_1 \cup \mathbf{M}_2 \setminus \{T_1, T_2\}$ with $|\mathbf{S}| = s \geq 0$. By Remark 1, we assume $T_1 \in \mathbf{M}_2$ and $T_2 \in \mathbf{M}_1$. Recall that the number of all variables is p . It follows that:

- In view of $|\mathbf{N} \cup \mathbf{S}| = m_1 + m_2 + s - 2 \triangleq m$, IAMB takes time $O[(mp+m)t_{n,m}]$ to finish an execution. Thus, the complexity of IAMB is $O(mpt_{n,m})$. KIAMB has almost the same complexity.
- For MIAMB, it first takes time $O[(m_1p + m_1)t_{n,m_1} + (m_2p + m_2)t_{n,m_2}]$ to find \mathbf{M}_1 and \mathbf{M}_2 ; then it seeks \mathbf{S} and refines \mathbf{N} taking time $O\{[s(p - m_1 - m_2 + 2) + m]t_{n,m}\}$. Hence, MIAMB needs time $O\{(m_1p + m_1)t_{n,m_1} + (m_2p + m_2)t_{n,m_2} + [s(p - m_1 - m_2 + 2) + m]t_{n,m}\}$ to finish an execution, so its complexity is $O(m_1pt_{n,m_1} + m_2pt_{n,m_2} + spt_{n,m})$. MKIAMB has almost the same complexity.

By this analysis, the complexity of MIAMB or MKIAMB is lower than that of IAMB or KIAMB. In fact, noting $t_{n,q}$ is exponential to q ($\leq m$; meaning $t_{n,q} \ll t_{n,m}$ in most situations) for $q = m_1, m_2$, this implies MIAMB/MKIAMB are expected to need much less time to run than IAMB/KIAMB, especially when \mathbf{T} contains many variables. The evaluation section (Figure 15) confirms this expectation in the case of moderately large sample size.

For the general case, using the notations in Subsection 4.2 with $|\mathbf{M}_i| = m_i$ ($i = 1, \dots, k$), we assume \mathbf{S}_i be an MBS for \mathbf{M}_{i-1}^* and \mathbf{M}_i , with $|\mathbf{S}_i| = s_i$ ($i = 2, \dots, k$). Denote $m_i^* \triangleq \sum_{j=1}^i m_j + \sum_{j=2}^i s_j - i$. Note that, in general, $t_{n,m_a} \ll t_{n,m_a^*} \ll t_{n,m_b^*}$ for $a < b$. Then, the IAMB or KIAMB algorithm has the complexity $O(m_k^*pt_{n,m_k^*})$, while MIAMB or MKIAMB has a lower complexity $O(\sum_{i=1}^k m_ipt_{n,m_i} + \sum_{i=2}^k s_ipt_{n,m_i^*})$.

According to this theoretical result on complexities, we can use the ordering, $\tau \triangleq \{i_1 \cdots, i_k\}$, in MIAMB or MKIAMB such that $m_{i_1} \leq \cdots \leq m_{i_k}$. This can reduce the complexities to some extent.

Besides, the additivity based MIAMB/MKIAMB algorithms have almost the same complexity as the MBS based MIAMB/MKIAMB, while the dummy MIAMB/MKIAMB have the complexity $O(m r_T p t_{n,m})$, where $m = \sum_{j=1}^k m_j$, $r_T = \prod_{j=1}^k r_{T_j}$, r_ξ denotes the number of configurations for ξ . It will be seen from Section 6 that, although the dummy MIAMB/MKIAMB are of high complexity theoretically, they usually perform well in multi-class prediction problems.

5. Benchmarking Study

This section makes a benchmarking study based on the data sets of six synthetic BNs. These data sets, generated by Tsamardinos et al. (2006) and Aliferis et al. (2010a), and the BNs are briefly described in Table 1. As Tsamardinos et al. (2006) and Aliferis et al. (2010a) stated, these BNs are representatives of a wide range of problem domains. Also, these BNs have different complexities (according to the number of nodes, the number of edges, maximal in-degree, maximal out-degree, and domain range). More details about the BNs and the used data sets are provided by Tsamardinos et al. (2006) and Aliferis et al. (2010a).

The following items are clarified before presenting the experimental results:

- *Measurements*: The primary measurement for the performance of an MB discovery algorithm used in our experiment is the weighted accuracy (WA), which is the average of the rate of true members and that of true nonmembers of an MB with respect to the truth. We also compute what we call the weighted precision (WP) as the average of the rate of true members and that of true nonmembers of an MB with respect to the output. In addition, we record the running time (RT) for every data set of each algorithm and for each BN. Here, RT refers to the single CPU time implemented on an Intel i7-3612QM 2.1 GHz and Windows 7 with 64 bits.

BN	Num. Nodes	Num. Edges	Maximal In-degree	Maximal Out-degree	Domain Range	Selected Targets	Sizes of Data Sets	Total RT (Hours)
Child10	200	257	2	7	2 ~ 6	$X_{131}, X_{132}, X_{98}, X_{194}, X_{184}, X_{22}, X_{135}, X_{60}$	5×500	3.4297
							1×5000	8.2220
ALARM10	370	570	4	7	2 ~ 4	$X_{341}, X_{48}, X_{37}, X_{249}, X_{209}, X_{188}, X_{192}, X_{161}$	5×500	4.6547
							1×5000	6.3721
Pigs	441	592	2	39	3 ~ 3	$X_{390}, X_{357}, X_{180}, X_{400}, X_{199}, X_{241}, X_{228}, X_{176}$	5×500	11.7651
							1×5000	15.8443
Link	724	1125	3	14	2 ~ 4	$X_{369}, X_{293}, X_{303}, X_{457}, X_{399}, X_{512}, X_{183}, X_{501}$	5×500	9.7932
							1×5000	16.1046
Lung Cancer	800	1476	4	28	2 ~ 3	$X_1, X_{416}, X_{345}, X_{641}, X_{513}, X_{198}, X_{78}, X_{746}$	5×500	16.0301
							1×5000	16.3790
Gene	801	977	4	10	3 ~ 5	$X_{801}, X_{301}, X_{569}, X_{317}, X_{185}, X_{622}, X_{516}, X_{577}$	5×500	17.2465
							1×5000	35.6790

Table 1: BNs and data sets.

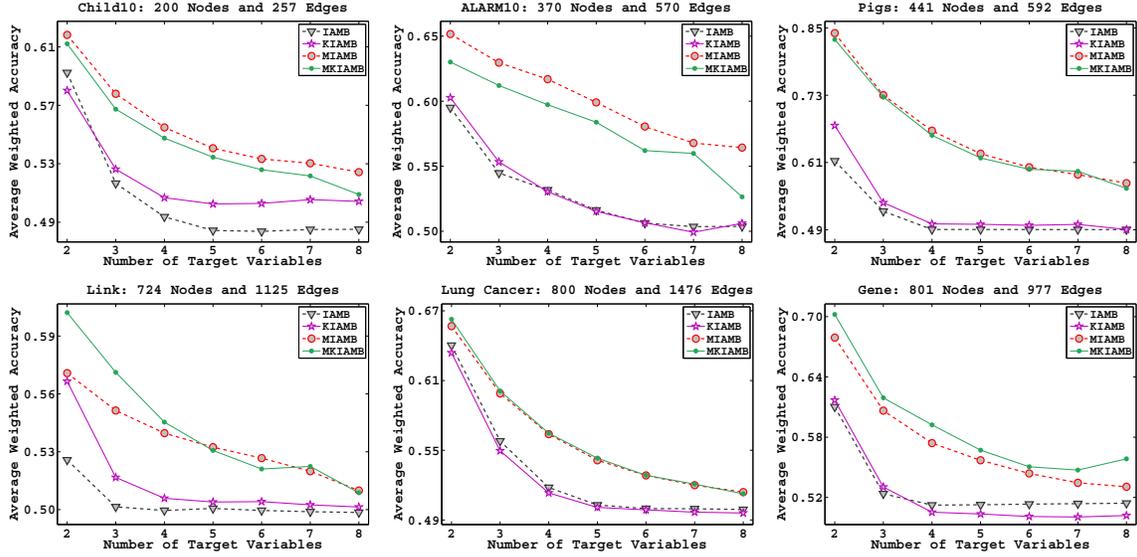


Figure 10: Average WA of the algorithms versus $|T|$ with respect to the data sets of size 500

- Algorithms:** Four algorithms are used: IAMB, KIAMB, MIAMB, and MKIAMB. We take $K = 0.8$ as the randomization parameter in KIAMB and MKIAMB due to the following two reasons: (i) Peña et al. (2007, p. 227) asserted that $K \in [0.7, 0.9]$ performs best; and (ii) $K = 0.8$ is an appropriate tradeoff between WA (or WP) and RT.
- Used CITest:** Following the *practical operation* suggested in Subsection 4.3, we implemented the algorithms via the G^2 test, and found G^2 is suitable for small $|T|$ but is not very suitable and even no longer works for large $|T|$. Then, we used the experimental CMI-based test with a relatively rough $\varepsilon \approx \sqrt{|T|} \cdot \varepsilon_0$, in which $\varepsilon_0 = 0.05$ if $n = 500$ and $\varepsilon_0 = 0.01$ if $n = 5000$; after that, we used the damped G^2 test by setting $\kappa = 5$. The results indicate both alternatives are desirable. Considering the association function, $f_D^{(2)}$ defined in (4) corresponding to the CMI-based test, contains no the average number, v , of values that each variable takes, we may need to reselect ε_0 for a BN with a very different v . For these reasons, we eventually decided to use the damped G^2 test for the four algorithms in our experiment.
- Data:** We use the data sets of sizes 500 and 5000, generated by Tsamardinos et al. (2006) and Aliferis et al. (2010a), which are available at <http://www.nyuinformatics.org/downloads/supplements/JMLR2009/index.html>.
- Targets:** We employ eight of those variables selected by Aliferis et al. (2010a, p. 226) as the potential targets for each BN. See Table 1 for details. Then, T is any possible combination of k targets for $k = 2, \dots, 8$.
- Steps:** For each BN with eight selected targets, the steps of making simulation based on the data set of size 5000 are as follows: (a) for $k = 2, \dots, 8$, call the four algorithms to obtain four MBs of $T \triangleq \{T_{i_1}, \dots, T_{i_k}\}$; (b) compute their WAs and WPs, and record the respective RTs; (c) take the average values of these $\binom{8}{k}$ WAs or WPs or RTs for each of the four algorithms. For the five data sets of size 500, each reported WA or WP or RT is the average value of the corresponding five results of an algorithm derived by (a) ~ (c) above.

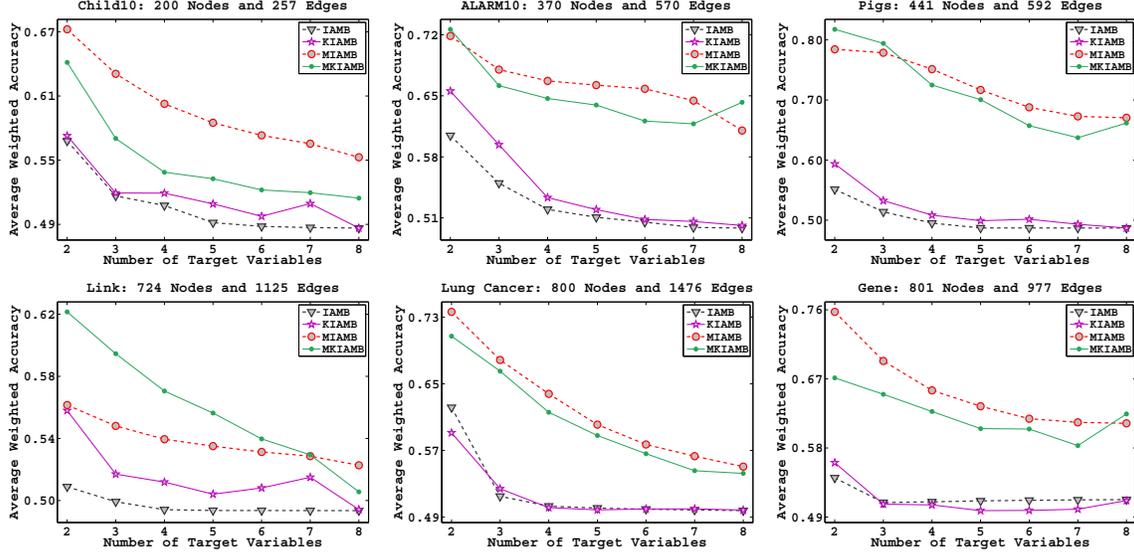


Figure 11: Average WA of the algorithms versus $|T|$ with respect to the data sets of size 5000

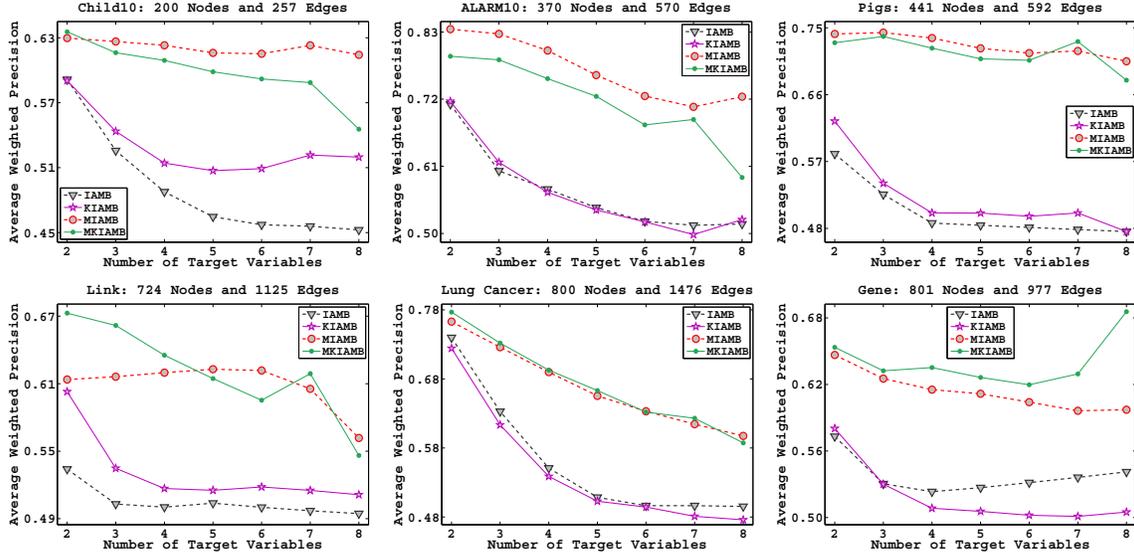


Figure 12: Average WP of the algorithms versus $|T|$ with respect to the data sets of size 500

According to the above description, we make computations with the aid of FullBNT (Murphy, 2007) and MIToolbox (Brown et al., 2012). The results of the WAs are presented in Figure 10 and Figure 11; the results of the WPs are given in Figure 12 and Figure 13; and the results of the RTs are shown in Figure 14, and Figure 15. The total RTs are presented in Table 1. By these figures, it is concluded that, on the whole, our MIAMB and MKIAMB have higher computational accuracies and lower time complexities than the existing IAMB and KIAMB.

Specifically, we have:

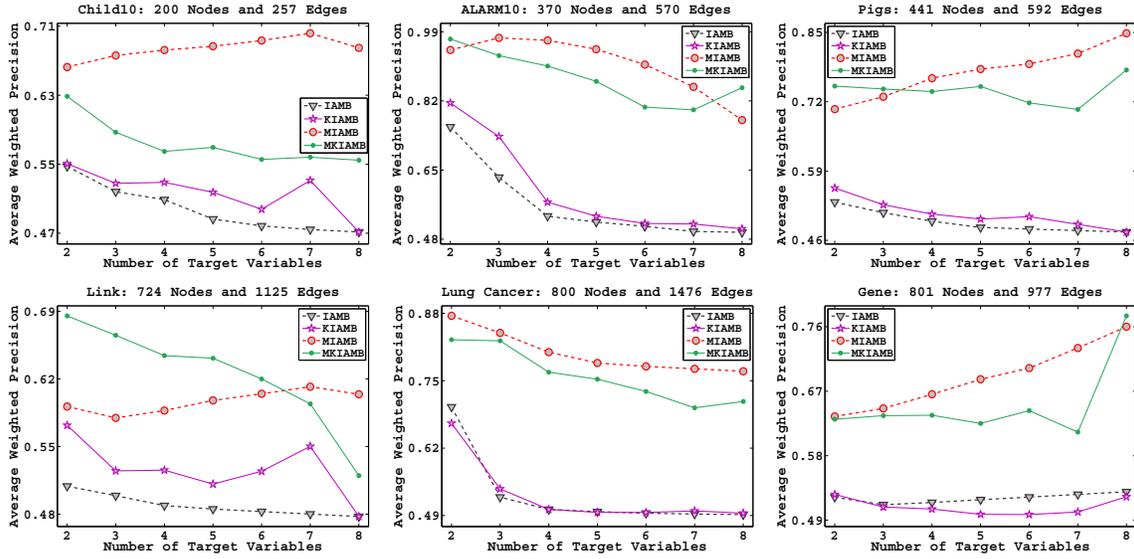


Figure 13: Average WP of the algorithms versus $|T|$ with respect to the data sets of size 5000

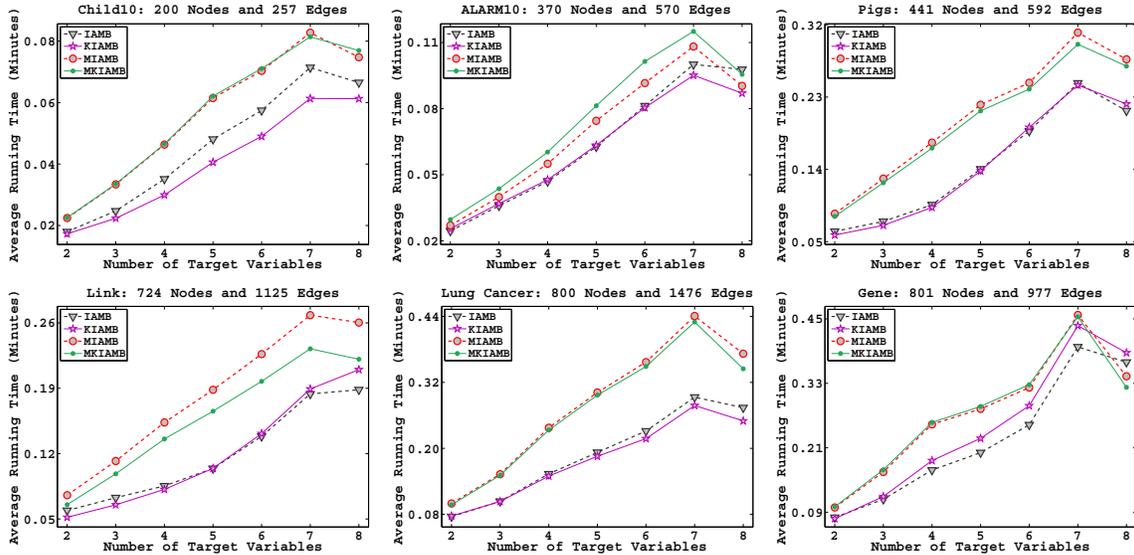


Figure 14: Average RT of the algorithms versus $|T|$ with respect to the data sets of size 500

- (i) *Performance on WA:* (a) MIAMB and MKIAMB have larger WAs than IAMB and KIAMB for all the six BNs in any case of $|T|$; (b) when $|T|$ increases, WA declines quickly for IAMB and KIAMB, but it decreases gently for MIAMB and MKIAMB; and (c) the improvements of MIAMB and MKIAMB over IAMB and KIAMB tend to be gradually noticeable and then reduce slightly as $|T|$ increases. The performance degradation along with the increase of $|T|$ can be attributed to two possible aspects: one is that the local composition assumption may be more apt to be violated for a larger $|T|$ because of synergy effects; and the other is that the assumption about

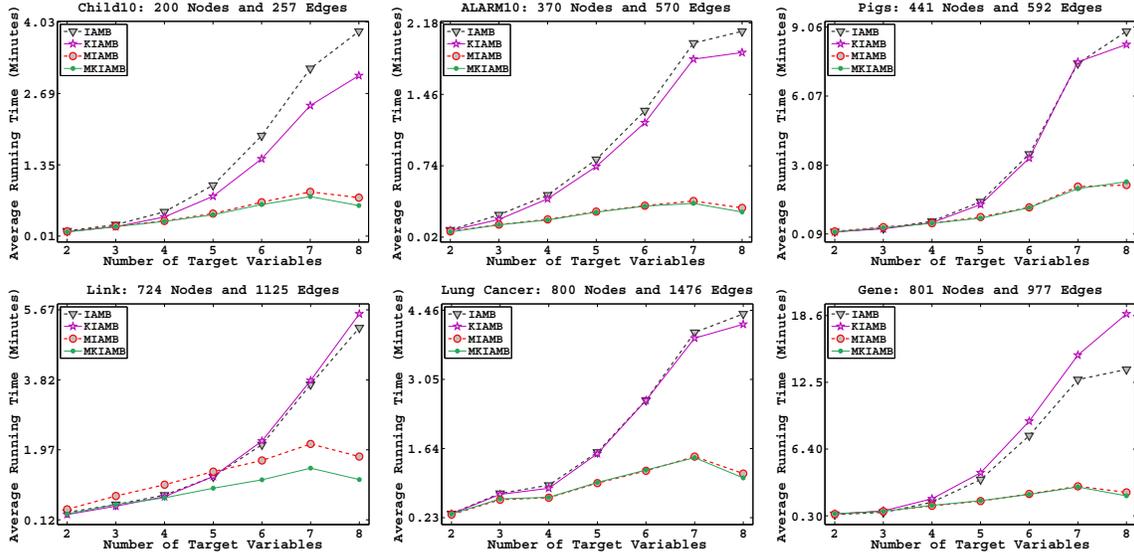


Figure 15: Average RT of the algorithms versus $|T|$ with respect to the data sets of size 5000

the correctness of CI tests may also be more apt to be violated for a larger $|T|$, due to the accumulation and propagation of the cascading errors (Bromberg and Margaritis, 2009). It is mentioned that (a)(b)(c) appear more evidently for the case of $n = 5000$ than for the case of $n = 500$.

- (ii) *Performance on WP*: The similar interpretations to (a)(c) of (i) are valid.
- (iii) *Performance on RT*: Here, we note that the real RT of an MB discovery algorithm is composed of two parts, in which the *part (I)* is for CI tests, and the *part (II)* is for all other computations. The part (I) is the major part used to measure the complexity of the MB discovery algorithm. Note also that the RT, $t_{n,q}$, of per CI test is linear to the sample size, n , and exponential to the number, q , of variables in the conditional set (see Subsection 4.4 for details).

This means that the part (II) of the real RT may dominate the part (I) if n is not large (for example, $n = 500$).

Let us now observe Figure 14 and Figure 15. First, both figures show the real RT that each algorithm needs is increasing along with the increase of $|T|$. Also, Figure 14 indicates MIAMB and MKIAMB need slightly longer time to run than IAMB and KIAMB, because the running for CI tests is dominated by the running for all other computations in the case of a small sample size, while Figure 15 reveals that the real RTs of IAMB and KIAMB increase sharply as $|T|$ increases and that the real RTs of MIAMB and MKIAMB increase slowly, just like the theoretical analyses about the complexities of the four algorithms show in Subsection 4.4.

In summary, the existing MB discovery algorithms, IAMB and KIAMB, can be approximately applied to the problem of MB discovery for multiple target variables when $|T|$ is small, but they will perform poorly if $|T|$ is moderately large. In comparison, our MIAMB and MKIAMB have higher accuracies and lower complexities for this problem, especially when $|T|$ is large.

6. Application to FS in Multi-Class Prediction Problems

In this section, we apply the MB discovery for multiple targets to FS in multi-class prediction problems based on a real data set, HIVA. This data set is very challenging in WCCI 2006 (<http://www.modelselect.inf.ethz.ch>) and IJCNN 2007 (<http://www.agnostic.inf.ethz.ch>), because it contains many very unbalanced variables.

Let $T \in V$ be a target variable taking values $\{1, \dots, t\}$, $t \geq 3$. The multi-class prediction problem is to select features of T from $V \setminus \{T\}$ such that T can be predicted as accurately as possible based on the chosen features. Let $\mathbf{T}^{(d)} \triangleq (T_1^{(d)}, \dots, T_t^{(d)})$ be the dummy version of T . Theoretically, $\mathbf{T}^{(d)}$ and T have the same MBs.

With these notations, the experiment is designed as follows:

- *Data*: HIVA contains 4229 data points and 1618 variables.
- *Targets*: In view of the fact that almost all variables in HIVA are binary, we randomly take k 2-class variables ($k = 2, \dots, 5$) to create a merged 2^k -class target, T . Accordingly, we rearrange the original data to get a data set that is used only for FS of T . Repeat this step $n \triangleq 200$ times. Denote the resulting targets and their dummy versions by T_1, \dots, T_n and $\mathbf{T}_1^{(d)}, \dots, \mathbf{T}_n^{(d)}$, respectively.
- *Algorithms*: We use the following 10 MB discovery algorithms to get the features for each T_j or $\mathbf{T}_j^{(d)} \triangleq (T_{j1}^{(d)}, \dots, T_{jt_k}^{(d)})$ with $t_k \triangleq 2^k$, $j = 1, \dots, n$:
 - a) IAMB/KIAMB-I: the IAMB/KIAMB algorithms working on T_j directly;
 - b) IAMB/KIAMB-II: the IAMB/KIAMB algorithms working on $\mathbf{T}_j^{(d)}$ (that is, with $\mathbf{T}_j^{(d)}$ as its multiple targets);
 - c) MIAMB/MKIAMB-I: the additivity based MIAMB/MKIAMB algorithms, taking the union of the outputs of IAMB/KIAMB with respect to $T_{ji}^{(d)}$ ($i = 1, \dots, t_k$) as its output;
 - d) MIAMB/MKIAMB-II: the MBS based MIAMB/MKIAMB algorithms, which are pseudo-coded in Algorithm 2;
 - e) MIAMB/MKIAMB-III: the dummy MIAMB/MKIAMB algorithms, which take the union of the outputs of IAMB/KIAMB with respect to $T_{ji}^{(d)}$ ($i = 1, \dots, t_k$) and removing redundant variables.
- *Classifier*: After making a number of preliminary experiments on the six benchmarking BNs, we found that the support vector machines (SVMs; implemented via LibSVM v3.22) perform the best in demonstrating the optimality of MBs for FS. This coincides with the assertion of Statnikov et al. (2013). Therefore, we use SVMs for our multi-class prediction problems. All the classifications are performed by 10-fold cross-validation.
- *Measurement of an algorithm*: For each target, the predictive quality of an MB is measured by the *balanced accuracy* defined as $\tau \triangleq \frac{1}{t_k} \sum_{\ell=1}^{t_k} (c_{\ell\ell} / \sum_{i=1}^{t_k} c_{i\ell})$, where $\mathbf{C} \triangleq (c_{i\ell})$ denotes the associated confusion matrix. As seen, τ is equal to one minus the *balanced error rate* used in WCCI 2006 and IJCNN 2007. We choose to use τ (instead of *ordinary accuracy*) because it trades off all values of the target in the sense that any unbalanced value (that the target

Problem	IAMB		MIAMB		
	I	II	I	II	III
4-class	0.9295 ± 0.082	0.9020 ± 0.113	0.9414 ± 0.068	0.9366 ± 0.073	0.9507 ± 0.057
8-class	0.9016 ± 0.102	0.8666 ± 0.133	0.9237 ± 0.091	0.9277 ± 0.087	0.9348 ± 0.077
16-class	0.8878 ± 0.105	0.8461 ± 0.143	0.9131 ± 0.087	0.9167 ± 0.086	0.9256 ± 0.075
32-class	0.8683 ± 0.113	0.8179 ± 0.157	0.9118 ± 0.078	0.9139 ± 0.078	0.9242 ± 0.067

Table 2: Balanced accuracy of IAMB/MIAMB algorithms in the form of “(mean ± std)”.

Problem	KIAMB		MKIAMB		
	I	II	I	II	III
4-class	0.9283 ± 0.085	0.8972 ± 0.124	0.9281 ± 0.092	0.9444 ± 0.066	0.9501 ± 0.058
8-class	0.9007 ± 0.105	0.8631 ± 0.144	0.9245 ± 0.089	0.9280 ± 0.082	0.9340 ± 0.078
16-class	0.8886 ± 0.106	0.8494 ± 0.142	0.9159 ± 0.084	0.9168 ± 0.086	0.9263 ± 0.075
32-class	0.8687 ± 0.114	0.8241 ± 0.151	0.9111 ± 0.081	0.9151 ± 0.076	0.9239 ± 0.068

Table 3: Balanced accuracy of KIAMB/MKIAMB algorithms in the form of “(mean ± std)”.

takes) should not impact on the accuracy too much.¹ On the other hand, when two outputs of algorithms have the same total numbers of “true positives + true negatives”, the balanced accuracy can identify the output that prefers to protect the scarce class as the better one, while the ordinary accuracy cannot. Finally, we compute the mean and standard deviation (std) of the n values of balanced accuracy, denoting them in the form of “(mean ± std)”.

The experiment is then performed following the above procedures. Its results are summarized in Table 2 and Table 3. In these two tables, the backcolor indicates the performance of algorithms with black corresponding to the best while light blue to the worst. By the results, it can be seen that MIAMB/MKIAMB outperform IAMB/KIAMB in most situations. Specifically, we have:

- IAMB/KIAMB algorithms: IAMB/KIAMB-I are much more preferred than IAMB/KIAMB-II.
- MIAMB/MKIAMB algorithms: MKIAMB-I has almost equal performance to KIAMB-I in 4-class problems, and they performs slightly better than IAMB/KIAMB-I in 16- and 32-class problems;

1. For example, consider an unbalanced target T and its classification with the following two confusion matrices (the left is *extremely bad*, while the right is *very good*):

Test \ Truth	$T = 1$	$T = 2$
$T = 1$	948	49
$T = 2$	2	1

Test \ Truth	$T = 1$	$T = 2$
$T = 1$	899	0
$T = 2$	51	50

Then, we have: (a) for the left *bad* confusion matrix, the ordinary accuracy equals 94.90% (meaning it is impacted deeply by the unbalanced value 1 of T), while its balanced accuracy equals 50.89%; (b) for the right *good* confusion matrix, its ordinary accuracy also equals 94.90%, but its balanced accuracy equals 97.32%. This means balanced accuracy is more reasonable than ordinary accuracy to measure classification performance for a practical problem containing unbalanced variables (note that such problems may frequently occur in practice).

Null hypothesis (H_0)	Problem			
	4-class	8-class	16-class	32-class
MIAMB-I \leq IAMB-I	2.0349×10^{-4}	1.3225×10^{-7}	1.3816×10^{-10}	4.0634×10^{-19}
MIAMB-II \leq IAMB-I	1.4772×10^{-2}	9.3393×10^{-10}	4.0911×10^{-12}	8.2111×10^{-21}
MIAMB-III \leq IAMB-I	2.1276×10^{-10}	9.5876×10^{-16}	1.8800×10^{-20}	9.1061×10^{-30}
MKIAMB-I \leq KIAMB-I	0.6221	4.9365×10^{-9}	5.9117×10^{-12}	2.0651×10^{-19}
MKIAMB-I = KIAMB-I	0.7558	—	—	—
MKIAMB-II \leq KIAMB-I	2.2839×10^{-6}	5.9538×10^{-11}	5.4564×10^{-12}	2.7021×10^{-22}
MKIAMB-III \leq KIAMB-I	4.5857×10^{-10}	1.1261×10^{-15}	4.7802×10^{-20}	6.4206×10^{-30}

Table 4: p -values on paired t -test for comparison between MIAMB/MKIAMB and IAMB/KIAMB. Here, the notations are defined as follows: letting \mathcal{A}_1 and \mathcal{A}_2 be two algorithms and P be a problem, if \mathcal{A}_1 is better (in the sense of possessing higher accuracy) than \mathcal{A}_2 when used to solve P , we denote it by $\mathcal{A}_1 > \mathcal{A}_2$ (w.r.t. P); otherwise, we write it as $\mathcal{A}_1 \leq \mathcal{A}_2$. In addition, we use $\mathcal{A}_1 = \mathcal{A}_2$ to denote $\mathcal{A}_1 \leq \mathcal{A}_2$ and $\mathcal{A}_1 \geq \mathcal{A}_2$.

MIAMB/MKIAMB-II significantly improve IAMB/KIAMB and even MIAMB/MKIAMB-I in most cases (although MIAMB/MKIAMB-II have larger std values than MIAMB/MKIAMB-I in some cases, the differences are slight). MIAMB/MKIAMB-III perform the best in all situations, with the highest mean values and the smallest std values.

Further, for any two algorithms, denote their balanced accuracy values as n ($= 200$) paired data points. Then, we can compute the p -values of *paired t -test* of associated hypotheses for one algorithm to be better (in the sense of possessing higher accuracy) than the other. The results are presented in Table 4. This table quantificationally shows the statistical significance of how much MIAMB/MKIAMB improve IAMB/KIAMB: in most cases, the improvement is more and more significant as the classification complex increases.

- The performance of each algorithm degrades with the increase of classification complexity. However, the degenerations of MIAMB/MKIAMB are slower than that of IAMB/KIAMB.

To compare IAMB/KIAMB and MIAMB/MKIAMB detailedly, we take the results of IAMB/KIAMB-I and MIAMB/MKIAMB-III to make a further analysis. For the 4-class prediction problem, denote the results of IAMB-I and MIAMB-III by $\tau_i^{(\text{IAMB})}$ and $\tau_i^{(\text{MIAMB})}$ for $i = 1, \dots, n$, and draw them in (a) of Figure 16. Put

$$\begin{aligned} I_1 &= \{i \in \{1, \dots, n\} : \tau_i^{(\text{IAMB})} > \tau_i^{(\text{MIAMB})}\}, \\ I_2 &= \{i \in \{1, \dots, n\} : \tau_i^{(\text{IAMB})} = \tau_i^{(\text{MIAMB})}\}, \\ I_3 &= \{i \in \{1, \dots, n\} : \tau_i^{(\text{IAMB})} < \tau_i^{(\text{MIAMB})}\}. \end{aligned}$$

Draw the scatters of $\tau_i^{(\text{IAMB})}$ and $\tau_i^{(\text{MIAMB})}$ for $i \in I_j$ in (a) of Figure 16. In addition, the information about (mean \pm std) of IAMB-I vs that of MIAMB-III is annotated in each title. For other three K -class prediction problems ($K = 8, 16, 32$), repeat the above steps to get the scatters drawn in the other subplots of Figure 16. Similarly, Figure 17 draws the results of KIAMB-I versus MKIAMB-III.

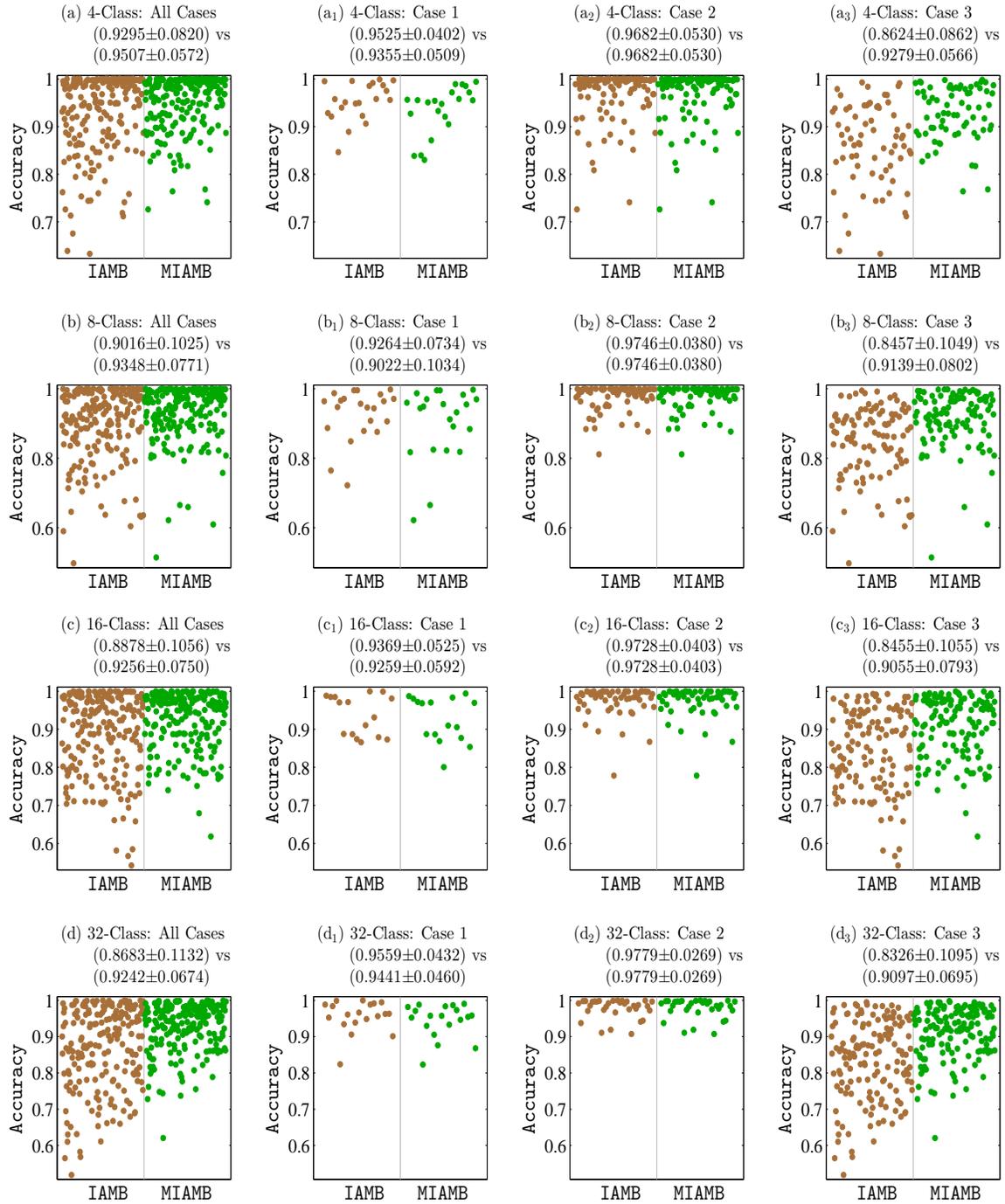


Figure 16: Balanced accuracy results on IAMB/MIAMB algorithms applied to 200 K -class prediction problems ($K = 4, 8, 16, 32$): the subplots in the first column for all the 200 results; the ones in the second column for the results that IAMB performs better than MIAMB; the ones in the third column for the results that IAMB and MIAMB perform equally well; the ones in the last column for the results that MIAMB performs better than IAMB.

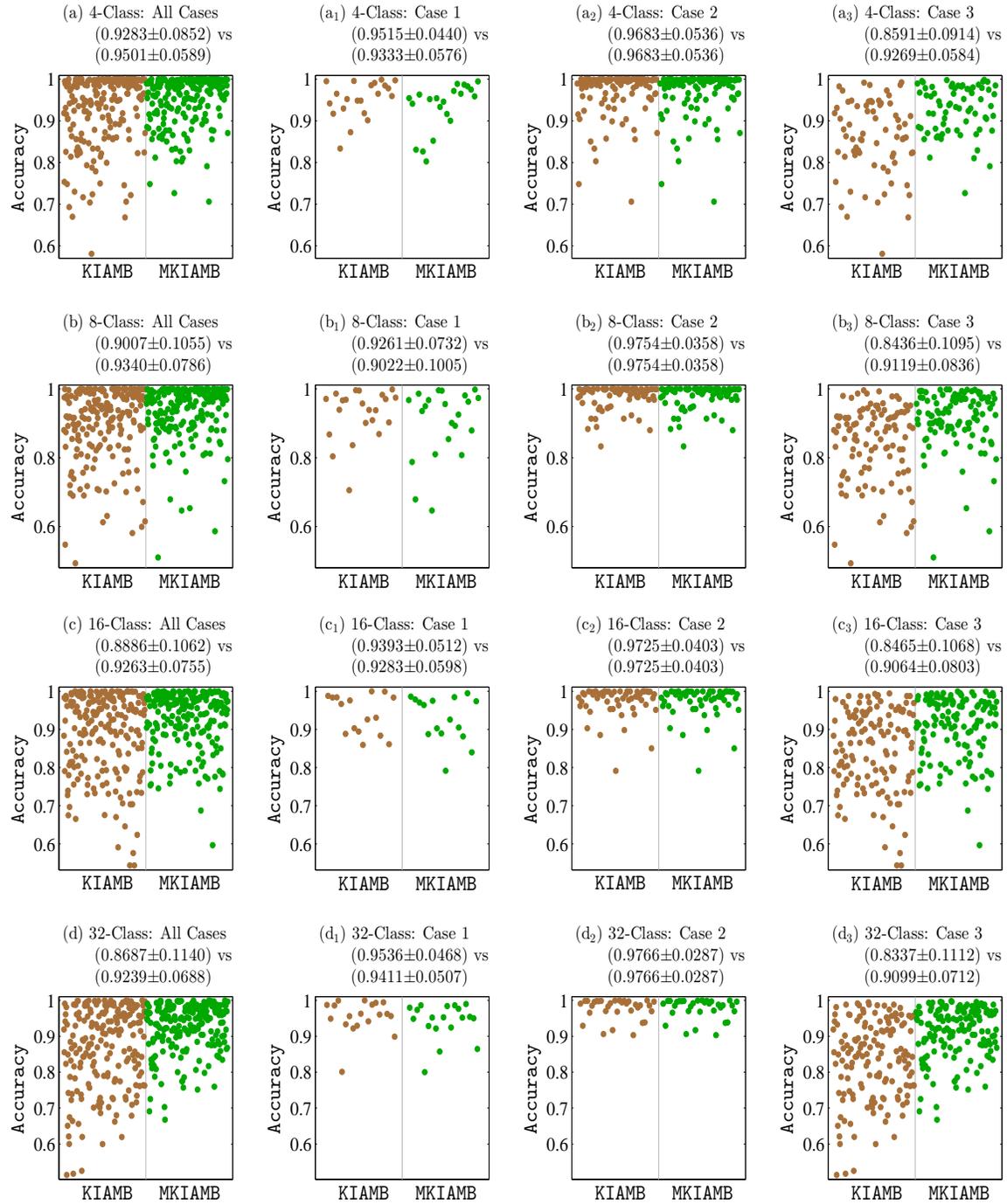


Figure 17: Balanced accuracy results on KIAMB/MKIAMB applied to 200 K -class prediction problems ($K = 4, 8, 16, 32$): the subplots in the first column for all the 200 results; the ones in the second column for the results that KIAMB performs better than MKIAMB; the ones in the third column for the results that KIAMB and MKIAMB perform equally well; the ones in the last column for the results that MKIAMB performs better than KIAMB.

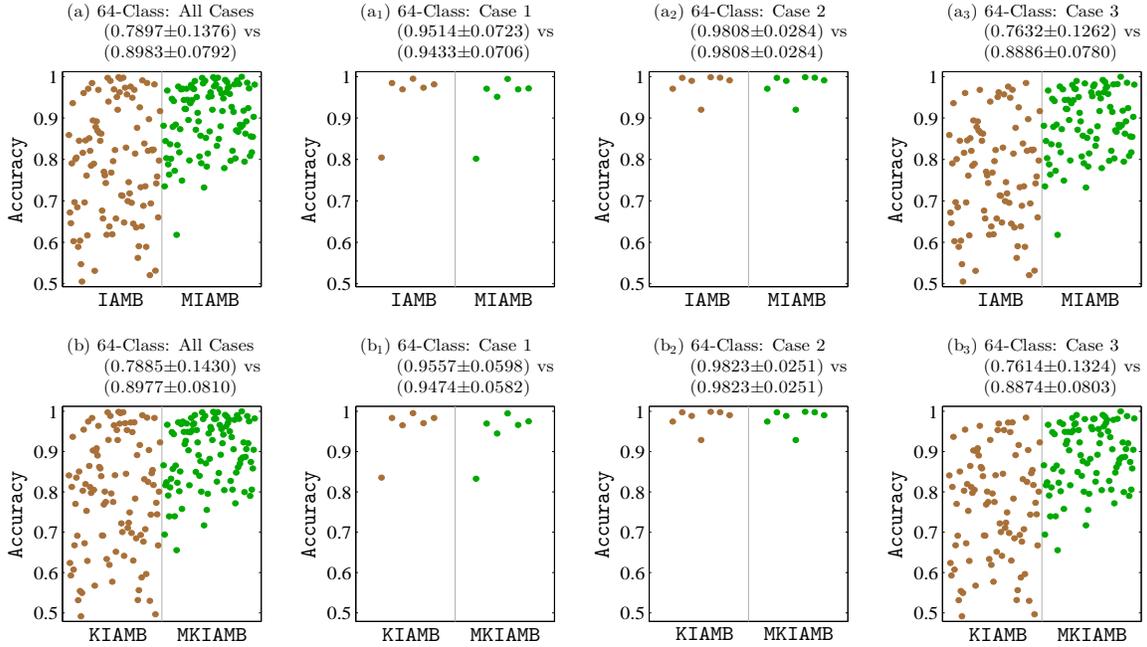


Figure 18: Balanced accuracy results on IAMB/MIAMB and KIAMB/MKIAMB applied to 100 64-class prediction problems: MIAMB/MKIAMB can improve IAMB/KIAMB substantially.

Figure 16 indicates that: (i) In most cases, MIAMB can improve IAMB to various degrees: MIAMB has higher mean values of balanced accuracy and smaller std values as well. Although there are a few of situations in which IAMB performs better than MIAMB, the difference of performance between them is very slight. In addition, there are some situations in which the two algorithms perform equally well (with very high mean and small std). (ii) MIAMB is more resistant to the classification complexity than IAMB: the improvements of MIAMB over IAMB become more and more visible with the increase of K (from 4 to 32). Figure 17 shows similar conclusions.

In brief, an FS problem for multi-class prediction can be transformed into a problem of MB discovery for multiple targets, and then get a more efficient solution. This idea may be particularly useful when the classification complexity is high or very high. To check this imagination, we apply the same procedures to 100 64-class prediction problems (also taken from the HIVA data set). The results are summarized in Figure 18, in which (a) and (a_{*j*}) are for IAMB/MIAMB while (b) and (b_{*j*}) are for KIAMB/MKIAMB, $j = 1, 2, 3$. By the figure, the improvement (nearly 14% on accuracy) of MIAMB/MKIAMB algorithms over IAMB/KIAMB is really desirable in the case of high classification complexity.

Finally, we apply LibSVM and the random forest (RF) algorithm (Breiman, 2001) to the whole HIVA data without any FS, considering LibSVM is of high classification performance while RF is a state-of-the-art FS algorithm. The results can be served as a baseline to see why FS (or equivalently, MB discovery) is necessary for a complex classification problem. Recall that HIVA contains many unbalanced variables, which enhance the classification complexity. Figure 19 draws the 95% confidence bands of LibSVM and FS, respectively.

By the figure, it follows that:

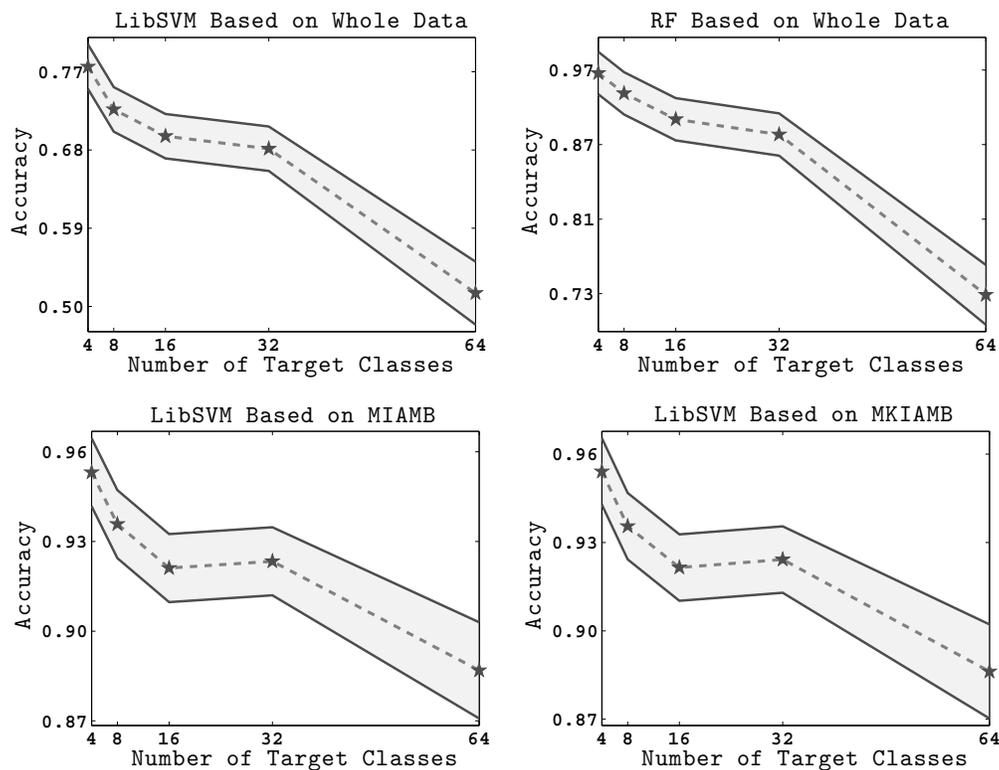


Figure 19: The 95% confidence bands of LibSVM and RF.

- Without any FS, LibSVM performs undesirably in all situations. This may be because too many noisy variables can lead to masking effects upon those unbalanced features such that LibSVM cannot classify targets expectedly. This shows the necessity of FS. In other words, LibSVM may not be suitable for some *high-dimension* problems, especially when there are many unbalanced variables.
- Without any FS, RF performs quite well when the classification complexity is not very high. However, with the increase of classification complexity, the performance of RF decreases gradually and then sharply. In other words, RF may not be suitable for those problems with too *high complexity*, especially when there are many unbalanced variables.

To observe why this happens, we check the results and then randomly take some targets (with extraordinarily low accuracy) to implement LibSVM and RF again by appropriately adjusting the algorithmic setting of LibSVM and increasing the number of trees of RF from 100 to 1000. However, the results change very little.

To make an intuitive comparison, the 95% confidence bands of MIAMB/MKIAMB-based LibSVM are also drawn in Figure 19. As seen, all methods degenerate with the increase of the classification complexity, but our methods degenerate far slower than LibSVM/RF based on the whole data. In a word, MB discovery (or equivalently, FS) is important to make classification, especially when the problem is of *high dimension* and of *high complexity*.

7. Concluding Remarks

In this paper, we considered the problem of Mb and MB of multiple variables. We first addressed their additivity under the local intersection assumption, and then studied this problem in the general case. The two algorithms that we proposed, MIAMB and MKIAMB, were proven to be correct under the local composition assumption with respect to single targets. The benchmarking study based on six synthetic BNs showed that MIAMB and MKIAMB have higher accuracies and lower complexities than the existing IAMB and KIAMB.

Before ending this paper, we present four concluding remarks as follows:

- (i) The first remark concerns a method of using MIAMB and MKIAMB to find an MB for a single variable. Such an idea is motivated by the following two aspects: (a) the local composition assumption may be violated in practice, and if this is the case, IAMB and KIAMB may perform not very well in MB discovery even for a single variable; (b) randomness of a data set may result in a violation to the assumption that all the CI tests involved are correct. Naturally, it is useful to take a remedy for these two situations. One remedial strategy is described as follows: letting $T \in V$ be the target variable, and M_1 be a potential MB discovered by IAMB or KIAMB, take $T_0 \triangleq \arg \max_{X \in M_1} f_D^{(\ell)}(T; X | M_1 \setminus \{X\})$ as a co-target of T for $\ell = 1$ or 2 or 3; then, employ MIAMB or MKIAMB to find a potential MB for $\{T, T_0\}$, saying M_2 . Finally, refine $\{T_0\} \cup M_2$ to obtain M , by virtue of the shrinking phase of IAMB or KIAMB, since this phase needs no the local composition precondition.
- (ii) Our MIAMB and MKIAMB contain an ordering τ , which may affect the RT and even the WA or WP. A question arises here: is there an optimal selection of τ such that MIAMB or MKIAMB has the highest accuracy and the lowest complexity?
- (iii) All the considered algorithms (IAMB, KIAMB, MIAMB, and MKIAMB) need the local composition assumption to theoretically guarantee their correctness. However, this precondition may be violated in practice and in this case only an approximate MB can be obtained by means of one of the above algorithms. Subsection 4.1 provides a potential remedy. We note that MIAMB and MKIAMB transform the problem of MB discovery for multiple targets into the ones for single targets. This idea provides a facilitation to use some stochastic optimization methods such as the particle swarm optimization algorithm (Kennedy and Eberhart, 1995, 1997).
- (iv) In Subsection 3.3, we provided a method for MB discovery of a *complex* single variable based on an MB discovery of some *simple* multiple variables. Let us now explain why this transformation method is efficient. With the notations used in Subsection 3.3, let

$$MB_T = MB_{T^{(d)}} \triangleq M.$$

Then, a variable X can enter and stay in M (in the sense of MB_T) if $T \not\perp\!\!\!\perp X | M \setminus \{X\}$. On the other hand, by Theorem 5, X can enter and stay in M (in the sense of $MB_{T^{(d)}}$) only if $T_j^{(d)} \not\perp\!\!\!\perp X | M \setminus \{X\}$ holds for some j . That is, we need to test the following two pairs of hypotheses:

$$\begin{aligned} H_0^{(1)} : T \perp\!\!\!\perp X | M \setminus \{X\} &\leftrightarrow H_1^{(1)} : T \not\perp\!\!\!\perp X | M \setminus \{X\}; \\ H_0^{(2)} : T_j^{(d)} \perp\!\!\!\perp X | M \setminus \{X\} &\leftrightarrow H_1^{(2)} : T_j^{(d)} \not\perp\!\!\!\perp X | M \setminus \{X\}. \end{aligned}$$

Clearly, when T is high-dimensional, the test for $H_0^{(1)} \leftrightarrow H_1^{(1)}$ requires far more data points than that for $H_0^{(2)} \leftrightarrow H_1^{(2)}$, since the test statistic for the first pair of hypotheses contains far more free parameters than that for the second. In addition, the transformation from T to $T^{(d)}$ can be easily made, with almost no running time. This explains why Theorem 5 is useful.

Acknowledgments

The authors are very grateful to the four anonymous reviewers and Prof. Marina Meila and Prof. Joris Mooij for their valuable comments and constructive suggestions which result in the present version. Thanks also to Prof. Kevin Murphy for all of his kind help.

This work was supported by the National Natural Science Foundation of China (61374183, 51472117, 51535005, 51675212), the Research Fund of State Key Laboratory of Mechanics and Control of Mechanical Structures (MCMS-0417G02, MCMS-0417G03), the Fundamental Research Funds for the Central Universities (NP2017101, NC2018001), the Project Funded by the Priority Academic Program Development of Jiangsu Higher Education Institutions, and the Open Fund for the Key Laboratory for Traffic and Transportation Security of Jiangsu Province.

Appendix A. Pseudo Codes for IAMB and KIAMB

This appendix presents the pseudo codes for IAMB and KIAMB. We mention here that, in Line 4 of KIAMB, M_2 denotes a random subset of M_1 with size $|M_2| = \max\{1, \lfloor |M_1| \cdot K \rfloor\}$.

Algorithm 3: IAMB and KIAMB	
<p>Procedure: $M \leftarrow \text{IAMB}(D; T)$ Input: D is a data matrix; T is a set of target variables. Output: an MB, M, of T.</p> <p>//Forward: Growing Phase 1 $M \leftarrow \emptyset$ 2 while M has changed do 3 $Y \leftarrow \arg \max_{X \in V \setminus M \setminus T} f_D(T; X M)$ 4 if $T \not\perp Y M$ then 5 $M \leftarrow M \cup \{Y\}$ 6 end 7 end</p> <p>//Backward: Shrinking Phase 8 foreach $X \in M$ do 9 if $T \perp X M \setminus \{X\}$ then 10 $M \leftarrow M \setminus \{X\}$ 11 end 12 end 13 return M</p>	<p>Procedure: $M \leftarrow \text{KIAMB}(D; T; K)$ Input: Besides $\{D, T\}$ as in IAMB, $K \in [0, 1]$ is a randomization parameter. Output: an Mb, M, of T.</p> <p>//Forward: Growing Phase 1 $M \leftarrow \emptyset$ 2 while M has changed do 3 if $M_1 \leftarrow \{X \in V \setminus M \setminus T : T \not\perp X M\} \neq \emptyset$ then 4 $Y \leftarrow \arg \max_{X \in M_2} f_D(T; X M)$ 5 $M \leftarrow M \cup \{Y\}$ 6 end 7 end</p> <p>//Backward: Shrinking Phase 8 foreach $X \in M$ do 9 if $T \perp X M \setminus \{X\}$ then 10 $M \leftarrow M \setminus \{X\}$ 11 end 12 end 13 return M</p>

Appendix B. Proofs

In this appendix, we give the proofs of the theoretical results.

Lemma 1 *The intersection property holds if and only if no information equivalence occurs.*

Proof Equivalently, we show that the intersection property is violated if and only if information equivalence occurs. The sufficiency holds clearly. To prove the necessity, we assume the intersection property is violated, that is, there are T , X , Y , and Z such that $T \perp\!\!\!\perp X|Z \cup Y$, and $T \perp\!\!\!\perp Y|Z \cup X$, but $T \not\perp\!\!\!\perp X \cup Y|Z$. Then, we can show $T \not\perp\!\!\!\perp X|Z$. In fact, if $T \perp\!\!\!\perp X|Z$, then combined with $T \perp\!\!\!\perp Y|Z \cup X$ and the contraction property, it concludes $T \perp\!\!\!\perp X \cup Y|Z$, which contradicts $T \not\perp\!\!\!\perp X \cup Y|Z$. Similarly, we can show $T \not\perp\!\!\!\perp Y|Z$. Therefore, X and Y are information equivalent with respect to T given Z . That is, information equivalence occurs. \blacksquare

Lemma 2 *For $T \subseteq V$, assume the type-II local condition holds. Then T has a unique MB.*

Proof Suppose T has two different MBs, M_1 and M_2 . Putting $M_{12} \triangleq M_1 \setminus M_2$, $M_{21} \triangleq M_2 \setminus M_1$, and $M \triangleq M_1 \cap M_2 \subsetneq M_i$ for $i = 1, 2$, we have

$$T \perp\!\!\!\perp V \setminus M_1 \setminus T | M_1 \Rightarrow T \perp\!\!\!\perp V \setminus M_1 \setminus T | M \cup M_{12}, \quad (8)$$

$$T \perp\!\!\!\perp V \setminus M_2 \setminus T | M_2 \Rightarrow T \perp\!\!\!\perp V \setminus M_2 \setminus T | M \cup M_{21}. \quad (9)$$

Now we show $T \not\perp\!\!\!\perp M_{12} | M$. In fact, suppose we have $T \perp\!\!\!\perp M_{12} | M$. This combined with (8) implies $T \perp\!\!\!\perp (V \setminus M_1 \setminus T) \cup M_{12} | M$, in view of the contraction property. Equivalently, $T \perp\!\!\!\perp V \setminus M \setminus T | M$, which contradicts the fact that M_1 is an MB of T , since $M \subsetneq M_1$. Hence, $T \not\perp\!\!\!\perp M_{12} | M$. Similarly, we can show $T \not\perp\!\!\!\perp M_{21} | M$. On the other hand, by the decomposition property, (9) and (8) indicate $T \perp\!\!\!\perp M_{12} | M \cup M_{21}$ and $T \perp\!\!\!\perp M_{21} | M \cup M_{12}$, respectively. Therefore, M_{12} and M_{21} are information equivalent with respect to T conditioned on M . This contradicts the precondition. The uniqueness of MB of T is shown under the type-II local condition. \blacksquare

Theorem 2 (Additivity of Mb) *Let (\mathbb{G}, \mathbb{P}) be a BN over V . The following two statements hold:*

- (i) *Let M_i be an Mb of $T_i \subseteq V$ for $i = 1, 2$, and assume $T_1 \cup T_2$ satisfies the local intersection assumption. Then, $(M_1 \cup M_2) \setminus (T_1 \cup T_2)$ is an Mb of $T_1 \cup T_2$.*
- (ii) *Let M_i be an Mb of $T_i \in V$ for $i = 1, \dots, k$, and assume $T \triangleq \{T_1, \dots, T_k\}$ satisfies the local intersection assumption. Then, $\bigcup_{i=1}^k M_i \setminus T$ is an Mb of T .*

Proof It suffices to prove (i), since (ii) is a direct consequence of (i) using induction on the number of variables involved in T .

Denote $V_0 = V \setminus N \setminus T$ with $N = N_1 \cup N_2$, in which $N_1 = M_1 \setminus T_2$ and $N_2 = M_2 \setminus T_1$. Note that N can also be expressed as $N = (M_1 \cup M_2) \setminus T$. First, we prove the following CI relationship, by means of the graphoid properties:

$$V_0 \perp\!\!\!\perp T_1 | N \cup T_2. \quad (10)$$

In fact, with the above notations, it is readily justified that $V \setminus M_1 \setminus T_1 = V_0 \cup [(N_2 \cup T_2) \setminus M_1]$. On the other hand, M_1 is an Mb of T_1 . Therefore, $T_1 \perp\!\!\!\perp V \setminus M_1 \setminus T_1 | M_1$, and thus we obtain $T_1 \perp\!\!\!\perp V_0 \cup [(N_2 \cup T_2) \setminus M_1] | M_1$. By the weak union property, $T_1 \perp\!\!\!\perp V_0 | M_1 \cup [(N_2 \cup T_2) \setminus M_1]$. This means (10) holds, since $M_1 \cup [(N_2 \cup T_2) \setminus M_1] = N \cup T_2$.

Similarly, $V_0 \perp\!\!\!\perp T_2 | N \cup T_1$, which combined with (10) indicates

$$V_0 \perp\!\!\!\perp T | N$$

by the local intersection assumption. Or equivalently, $T \perp\!\!\!\perp V \setminus N \setminus T | N$. That is, $(M_1 \cup M_2) \setminus T = N$ is an Mb of T . The proof is completed. \blacksquare

Remark 1 *In the case of either $T_1 \subseteq V \setminus M_2$ or $T_2 \subseteq V \setminus M_1$, the conclusion of (i) in Theorem 2 holds without requiring the local intersection assumption.*

Proof If $T_1 \subseteq V \setminus M_2$ but $T_2 \not\subseteq V \setminus M_1$, $M_1 \cup M_2$ is then an Mb of T_1 according to the weak union property whereas $(M_1 \setminus T_2) \cup M_2$ is an Mb of T_2 . Equivalently, we have

$$\begin{aligned} T_1 &\perp\!\!\!\perp V \setminus (M_1 \cup M_2) \setminus T_1 | M_1 \cup M_2, \\ T_2 &\perp\!\!\!\perp V \setminus [(M_1 \setminus T_2) \cup M_2] \setminus T_2 | (M_1 \setminus T_2) \cup M_2. \end{aligned}$$

By means of the contraction property and the decomposition property, it is seen that

$$\begin{aligned} V \setminus [(M_1 \cup M_2) \setminus T_2] \setminus T &\perp\!\!\!\perp T_1 | [(M_1 \cup M_2) \setminus T_2] \cup T_2, \\ V \setminus [(M_1 \cup M_2) \setminus T_2] \setminus T &\perp\!\!\!\perp T_2 | (M_1 \cup M_2) \setminus T_2, \end{aligned}$$

so $T \perp\!\!\!\perp V \setminus [(M_1 \cup M_2) \setminus T_2] \setminus T | (M_1 \cup M_2) \setminus T_2$. That is, $(M_1 \cup M_2) \setminus (T_1 \cup T_2) = (M_1 \cup M_2) \setminus T_2$ is an Mb of T . If $T_1 \not\subseteq V \setminus M_2$ but $T_2 \subseteq V \setminus M_1$, we can similarly show

$$(M_1 \cup M_2) \setminus (T_1 \cup T_2) = (M_1 \cup M_2) \setminus T_1$$

is an Mb of T . Finally, if $T_1 \subseteq V \setminus M_2$ and $T_2 \subseteq V \setminus M_1$, imposing decomposition on

$$T_1 \perp\!\!\!\perp V \setminus (M_1 \cup M_2) \setminus T_1 | M_1 \cup M_2$$

and weak union on $T_2 \perp\!\!\!\perp V \setminus (M_1 \cup M_2) \setminus T_2 | M_1 \cup M_2$, we get

$$\begin{aligned} V \setminus (M_1 \cup M_2) \setminus T &\perp\!\!\!\perp T_1 | M_1 \cup M_2, \\ V \setminus (M_1 \cup M_2) \setminus T &\perp\!\!\!\perp T_2 | (M_1 \cup M_2) \cup T_1. \end{aligned}$$

By the contraction property, $T \perp\!\!\!\perp V \setminus (M_1 \cup M_2) \setminus T | M_1 \cup M_2$. That is,

$$(M_1 \cup M_2) \setminus (T_1 \cup T_2) = M_1 \cup M_2$$

is an Mb of T . The conclusion is proved. \blacksquare

Theorem 3 (Additivity of MB) *Let (\mathbb{G}, \mathbb{P}) be a BN over V . The following two statements hold:*

- (i) *Assume $T_1 \cup T_2$ satisfies the local intersection assumption. Let M_i be the unique MB of T_i for $i = 1, 2$. Then, $(M_1 \cup M_2) \setminus (T_1 \cup T_2)$ is the unique MB of $T_1 \cup T_2$.*
- (ii) *Assume $T \triangleq \{T_1, \dots, T_k\}$ satisfies the local intersection assumption. Let M_i be the unique MB of T_i for $i = 1, \dots, k$. Then, $\bigcup_{i=1}^k M_i \setminus T$ is the unique MB of T .*

Proof We need only to prove (i), since (ii) is a direct consequence of (i).

Denote $N_1 = \mathbf{M}_1 \setminus \mathbf{T}_2$ and $N_2 = \mathbf{M}_2 \setminus \mathbf{T}_1$. By Theorem 2, $(\mathbf{M}_1 \cup \mathbf{M}_2) \setminus \mathbf{T} = N_1 \cup N_2$ is an Mb of \mathbf{T} . Therefore, it suffices to prove the minimality of $N_1 \cup N_2$, based on Lemma 2. In fact, let N_0 be any Mb of \mathbf{T} which is a subset of $N_1 \cup N_2$. Note that $\mathbf{T}_i \cap N_0 = \emptyset$ for $i = 1, 2$. Denote now $\mathbf{M} = N_0 \cup (\mathbf{M}_1 \setminus N_1) = N_0 \cup (\mathbf{M}_1 \cap \mathbf{T}_2)$. It follows that

- \mathbf{M}_1 is the MB of \mathbf{T}_1 : This implies $\mathbf{T}_1 \perp\!\!\!\perp V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1 \mid \mathbf{M}_1$, or equivalently, we have

$$\mathbf{T}_1 \perp\!\!\!\perp V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1 \mid (\mathbf{M}_1 \cap \mathbf{M}) \cup (\mathbf{M}_1 \setminus \mathbf{M}), \quad (11)$$

in view of $\mathbf{M}_1 = (\mathbf{M}_1 \cap \mathbf{M}) \cup (\mathbf{M}_1 \setminus \mathbf{M})$.

- N_0 is an Mb of \mathbf{T} : Equivalently, we have $\mathbf{T}_1 \cup \mathbf{T}_2 \perp\!\!\!\perp V \setminus N_0 \setminus \mathbf{T}_2 \setminus \mathbf{T}_1 \mid N_0$, which gives

$$\mathbf{T}_1 \perp\!\!\!\perp V \setminus (N_0 \cup \mathbf{T}_2) \setminus \mathbf{T}_1 \mid N_0 \cup \mathbf{T}_2,$$

according to the weak union property, and thus $\mathbf{T}_1 \perp\!\!\!\perp V \setminus (\mathbf{M} \cup \mathbf{T}_2) \setminus \mathbf{T}_1 \mid \mathbf{M} \cup \mathbf{T}_2$ in view of $N_0 \cup \mathbf{T}_2 = \mathbf{M} \cup \mathbf{T}_2$, or equivalently we have $\mathbf{T}_1 \perp\!\!\!\perp V \setminus \mathbf{M} \setminus \mathbf{T}_1 \setminus \mathbf{T}_2 \mid \mathbf{M} \cup \mathbf{T}_2$. By the self-conditioning property, this leads to $\mathbf{T}_1 \perp\!\!\!\perp V \setminus (\mathbf{M}_1 \cap \mathbf{M}) \setminus \mathbf{T}_1 \mid \mathbf{M} \cup \mathbf{T}_2$. Therefore,

$$\mathbf{T}_1 \perp\!\!\!\perp (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1) \cup (\mathbf{M}_1 \setminus \mathbf{M}) \mid \mathbf{M} \cup \mathbf{T}_2,$$

in terms of $V \setminus (\mathbf{M}_1 \cap \mathbf{M}) \setminus \mathbf{T}_1 = (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1) \cup (\mathbf{M}_1 \setminus \mathbf{M})$. By the weak union property, this indicates $\mathbf{T}_1 \perp\!\!\!\perp \mathbf{M}_1 \setminus \mathbf{M} \mid (\mathbf{M} \cup \mathbf{T}_2) \cup (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1)$. Consequently,

$$\mathbf{T}_1 \perp\!\!\!\perp \mathbf{M}_1 \setminus \mathbf{M} \mid (\mathbf{M}_1 \cap \mathbf{M}) \cup (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1), \quad (12)$$

due to $(\mathbf{M} \cup \mathbf{T}_2) \cup (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1) = (\mathbf{M}_1 \cap \mathbf{M}) \cup (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1)$.

By the local intersection property, (11)(12) indicate $\mathbf{T}_1 \perp\!\!\!\perp (\mathbf{M}_1 \setminus \mathbf{M}) \cup (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1) \mid \mathbf{M}_1 \cap \mathbf{M}$, so

$$\mathbf{T}_1 \perp\!\!\!\perp V \setminus (\mathbf{M}_1 \cap \mathbf{M}) \setminus \mathbf{T}_1 \mid \mathbf{M}_1 \cap \mathbf{M},$$

since $(\mathbf{M}_1 \setminus \mathbf{M}) \cup (V \setminus \mathbf{M}_1 \setminus \mathbf{T}_1) = V \setminus (\mathbf{M}_1 \cap \mathbf{M}) \setminus \mathbf{T}_1$. Hence, $\mathbf{M}_1 \cap \mathbf{M} (\subseteq \mathbf{M}_1)$ is an Mb of \mathbf{T}_1 . On the other hand, \mathbf{M}_1 is the MB of \mathbf{T}_1 and thereby $\mathbf{M}_1 \cap \mathbf{M} = \mathbf{M}_1$, or equivalently,

$$N_1 \cup (\mathbf{M}_1 \cap \mathbf{T}_2) = \mathbf{M}_1 \subseteq \mathbf{M} = N_0 \cup (\mathbf{M}_1 \cap \mathbf{T}_2),$$

which means $N_1 \subseteq N_0$. In a similar fashion, $N_2 \subseteq N_0$. Combined with $N_0 \subseteq N_1 \cup N_2$, the expected relationship $N_0 = N_1 \cup N_2$ follows. This indicates that $N_1 \cup N_2$ is an MB of \mathbf{T} . The proof is completed, since Lemma 2 shows the uniqueness of MB under the local intersection assumption. ■

Theorem 4 Assume $\mathbf{S} \subseteq V \setminus N \setminus \mathbf{T}$. Then, the following statements are equivalent:

- \mathbf{S} is an MbS;
- $\mathbb{I}(\mathbf{T}_1; \mathbf{T}_2 \mid N \cup \mathbf{S}) = \min_{\mathbf{S}' \subseteq V \setminus N \setminus \mathbf{T}} \mathbb{I}(\mathbf{T}_1; \mathbf{T}_2 \mid N \cup \mathbf{S}')$;
- $\mathbb{I}(\mathbf{T}; \mathbf{S} \mid N) = \max_{\mathbf{S}' \subseteq V \setminus N \setminus \mathbf{T}} \mathbb{I}(\mathbf{T}; \mathbf{S}' \mid N)$;
- $N \cup \mathbf{S}$ is an Mb of \mathbf{T}_1 in $V \setminus \mathbf{T}_2$ (or $N \cup \mathbf{S}$ is an Mb of \mathbf{T}_2 in $V \setminus \mathbf{T}_1$).

In addition, if S is an MbS, then it is also an MBS if and only if $T_1 \perp\!\!\!\perp Y | N \cup (S \setminus \{Y\})$ or $T_2 \perp\!\!\!\perp Y | N \cup (S \setminus \{Y\})$ holds for any $Y \in S$.

Proof We first prove (i) \Leftrightarrow (ii). Put $Q = V \setminus (N \cup S) \setminus T$. First, note that $T_1 \perp\!\!\!\perp V \setminus M_1 \setminus T_1 | M_1$ since M_1 is an Mb of T_1 . In other words, $a \triangleq \mathbb{I}(T_1; Q \cup (V \setminus M_1 \setminus T_1 \setminus Q) | M_1) = 0$. By the chain rule for CMI (Cover and Thomas, 2006), we have $\mathbb{I}(T_1; V \setminus M_1 \setminus T_1 \setminus Q | M_1) = 0$. It follows that

$$\begin{aligned}
 a &= \mathbb{I}(T_1; V \setminus M_1 \setminus T_1 \setminus Q | M_1) + \mathbb{I}(T_1; Q | (V \setminus M_1 \setminus T_1 \setminus Q) \cup M_1) \\
 &= \mathbb{I}(T_1; Q | (V \setminus M_1 \setminus T_1 \setminus Q) \cup M_1) \\
 &= \mathbb{I}(T_1; Q | N \cup S \cup T_2) \\
 &= \mathbb{I}(T; Q | N \cup S) - \mathbb{I}(T_2; Q | N \cup S) \\
 &\triangleq b - c,
 \end{aligned} \tag{13}$$

which combined with $a = 0$ gives $b = c$. Observing $T_2 \perp\!\!\!\perp V \setminus M_2 \setminus T_2 | M_2$ since M_2 is an Mb of T_2 , we obtain $T_2 \perp\!\!\!\perp Q | (V \setminus M_2 \setminus T_2 \setminus Q) \cup M_2$ by using the weak union property, or equivalently, $T_2 \perp\!\!\!\perp Q | N \cup S \cup T_1$, so $\mathbb{I}(T_2; Q | N \cup S \cup T_1) = 0$. This means

$$\begin{aligned}
 0 \leq c &= \mathbb{I}(T_2; Q | N \cup S) \\
 &= \mathbb{I}(T_2; T_1 \cup Q | N \cup S) - \mathbb{I}(T_2; T_1 | N \cup S \cup Q) \\
 &= \mathbb{I}(T_2; T_1 | N \cup S) + \mathbb{I}(T_2; Q | N \cup S \cup T_1) - \mathbb{I}(T_2; T_1 | N \cup S \cup Q) \\
 &= \mathbb{I}(T_2; T_1 | N \cup S) - \mathbb{I}(T_2; T_1 | N \cup S \cup Q).
 \end{aligned} \tag{14}$$

- (i) \Leftarrow (ii): If $\mathbb{I}(T_1; T_2 | N \cup S) \leq \mathbb{I}(T_1; T_2 | N \cup S')$ holds for any $S' \subseteq V \setminus N \setminus T$, then (14) indicates $c = 0$ since $0 \leq c = \mathbb{I}(T_2; T_1 | N \cup S) - \mathbb{I}(T_2; T_1 | N \cup S \cup Q) \leq 0$. Therefore,

$$\mathbb{I}(T; V \setminus (N \cup S) \setminus T | N \cup S) = \mathbb{I}(T; Q | N \cup S) = b = 0,$$

because of $b = c$. That is, $T \perp\!\!\!\perp V \setminus (N \cup S) \setminus T | N \cup S$, which means $N \cup S$ is an Mb of T , or equivalently, S is an MbS.

- (i) \Rightarrow (ii): Observe that $\mathbb{I}(T_1; T_2 | V \setminus T) = \min_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T_1; T_2 | N \cup S')$ holds according to (14) holding for any $S \subseteq V \setminus N \setminus T$ and $N \cup S \cup Q = V \setminus T$. Then,

$$\mathbb{I}(T; Q | N \cup S) = \mathbb{I}(T; V \setminus (N \cup S) \setminus T | N \cup S) = 0$$

follows immediately if $N \cup S$ is an Mb of T . By (13) and (14), we have

$$\mathbb{I}(T_1; T_2 | N \cup S) = \mathbb{I}(T_1; T_2 | V \setminus T) = \min_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T_1; T_2 | N \cup S'), \tag{15}$$

noting again $N \cup S \cup Q = V \setminus T$. This means $N \cup S$ is an MbS.

To prove the equivalence between (ii) and (iii), we need only to show

$$\mathbb{I}(T_1; T_2 | N \cup S) = \mathbb{I}(T_1; M_1) + \mathbb{I}(T_2; M_2) - \mathbb{I}(T_1; T_2) - \mathbb{I}(T; N \cup S). \tag{16}$$

In fact, using $\mathbb{I}(T; N \cup S) = \mathbb{I}(T_1; N \cup S) + \mathbb{I}(T_2; N \cup S | T_1)$, we have

$$\begin{aligned}
 d &\triangleq \mathbb{I}(T_1; T_2 | N \cup S) + \mathbb{I}(T_1; T_2) + \mathbb{I}(T; N \cup S) \\
 &= \mathbb{I}(T_1; T_2 | N \cup S) + \mathbb{I}(T_1; N \cup S) + \mathbb{I}(T_2; T_1) + \mathbb{I}(T_2; N \cup S | T_1) \\
 &= \mathbb{I}(T_1; T_2 \cup N \cup S) + \mathbb{I}(T_2; N \cup S \cup T_1) \\
 &= \mathbb{I}(T_1; M_1) + \mathbb{I}(T_2; M_2),
 \end{aligned}$$

which is equivalent to (16). This means (ii) \Leftrightarrow (iii).

Now, we show that (i) is equivalent to (iv):

- (i) \Rightarrow (iv): This implication holds clearly due to the decomposition property.
- (i) \Leftarrow (iv): Assume $N \cup S$ is an Mb of T_1 in $V \setminus T_2$, that is, $T_1 \perp (V \setminus T_2) \setminus (N \cup S) \mid N \cup S$, or equivalently, $V \setminus (N \cup S) \setminus T \perp T_1 \mid N \cup S$. On the other hand, M_2 is an Mb of T_2 in V , meaning $T_2 \perp V \setminus M_2 \setminus T_2 \mid M_2$, which combined with the weak union property gives $V \setminus (N \cup S) \setminus T \perp T_2 \mid N \cup S \cup T_1$. By the contraction property, $T \perp V \setminus (N \cup S) \setminus T \mid N \cup S$. This means S is an MbS. Similarly, if $N \cup S$ is an Mb of T_2 in $V \setminus T_1$, we can show S is an MbS.

Finally, we prove that an MbS, S , is an MBS if and only if $T_2 \not\perp Y \mid N \cup (S \setminus \{Y\})$ holds for any $Y \in S$. We first prove the necessity by reductio ad absurdum. Suppose there is some variable $Y \in S$ such that $T_2 \perp Y \mid N \cup R$, in which $R \triangleq S \setminus \{Y\}$. Recall that S is an MbS, we get

$$T_2 \perp (V \setminus T_1) \setminus (N \cup S) \setminus T_2 \mid N \cup S.$$

Equivalently, $T_2 \perp V \setminus (N \cup R) \setminus T \setminus \{Y\} \mid (N \cup R) \cup \{Y\}$, which combined with $T_2 \perp Y \mid N \cup R$ gives $T_2 \perp V \setminus (N \cup R) \setminus T \mid N \cup R$, in view of the contraction property. That is, $T_2 \perp (V \setminus T_1) \setminus (N \cup R) \setminus T_2 \mid N \cup R$. Therefore, $N \cup R$ is an Mb of T_2 in $V \setminus T_1$, and thus an MbS of T to N . This contradicts the condition that S is an MBS of T to N , and thus $T_2 \not\perp Y \mid N \cup (S \setminus \{Y\})$ holds for any $Y \in S$.

To prove the sufficiency, we suppose S is not a MBS of T to N , that is, there is some $R \subsetneq S$ such that R is an MbS of T to N . Take any given variable, Y , in $S \setminus R$. Then, $N \cup R$ is an Mb of T_2 in $V \setminus T_1$. That is, $T_2 \perp (V \setminus T_1) \setminus (N \cup R) \setminus T_2 \mid N \cup R$. By the weak union property, we have

$$T_2 \perp (V \setminus T_1) \setminus [N \cup (S \setminus \{Y\})] \setminus T_2 \mid N \cup (S \setminus \{Y\}).$$

This combined with the decomposition property means $T_2 \perp Y \mid N \cup (S \setminus \{Y\})$, since

$$Y \in (V \setminus T_1) \setminus [N \cup (S \setminus \{Y\})] \setminus T_2,$$

and thus leads to a contradiction to the condition that $T_2 \not\perp Y \mid N \cup (S \setminus \{Y\})$ holds for any $Y \in S$.

The proof of Theorem 4 is completed. \blacksquare

Example 1 Consider the BN (\mathbb{G}, \mathbb{P}) over $V = \{A, B, C, D\}$ presented in Figure 8, in which A, B , and C take $\{1, 2, 3\}$ while D takes $\{1, 2\}$. Put $T = \{T_1, T_2\}$, $N = (M_1 \cup M_2) \setminus T = \emptyset$, and $S = \{C\}$, $S_0 = \{C, D\}$ with $T_1 = A$, $T_2 = B$, $M_1 = \{B\}$, $M_2 = \{A\}$. By Figure 3, we can easily conclude that A and B are information equivalent with respect to C . This means

$$\mathbb{I}(C; A) > 0, \quad \mathbb{I}(C; B \mid A) = 0; \quad \text{and} \quad \mathbb{I}(C; B) > 0, \quad \mathbb{I}(C; A \mid B) = 0. \quad (17)$$

It follows from the chain rule for CMI (Cover and Thomas, 2006) that

(i) M_1 is an MB of T_1 in V : By (17), we have

- $\mathbb{I}(A; C, D \mid B) = \mathbb{I}(A; C \mid B) + \mathbb{I}(A; D \mid B, C) = 0$, since $\{B, C\}$ d -separates $\{A\}$ and $\{D\}$;
- $\mathbb{I}(A; C, D) \geq \mathbb{I}(A; C) > 0$.

(ii) M_2 is an MB of T_2 in V : By (17), we have

- $\mathbb{I}(B; C, D | A) = \mathbb{I}(B; C | A) + \mathbb{I}(B; D | A, C) = 0$, since $\{A, C\}$ d -separates $\{B\}$ and $\{D\}$;
- $\mathbb{I}(B; C, D) \geq \mathbb{I}(B; C) > 0$.

(iii) $N \cup S$ is an Mb of T in V , so S is an MbS: By (17),

$$\mathbb{I}(A, B; D | C) = \mathbb{I}(A; D | C) + \mathbb{I}(B; D | A, C) = 0,$$

because $\{C\}$ d -separates $\{A\}$ and $\{D\}$, while $\{A, C\}$ d -separates $\{B\}$ and $\{D\}$.

(iv) $\mathbb{I}(T_1; T_2 | N \cup S) = \min_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T_1; T_2 | N \cup S')$: it suffices to show the following inequalities:

- $\mathbb{I}(A; B | C) = \mathbb{I}(A; B | C, D)$. In fact,

$$\begin{aligned} \mathbb{I}(A; B | C, D) &= \mathbb{I}(A; B, D | C) - \mathbb{I}(A; D | C) \\ &= \mathbb{I}(A; B | C) + \mathbb{I}(A; D | B, C) - \mathbb{I}(A; D | C) = \mathbb{I}(A; B | C), \end{aligned}$$

since both $\{B, C\}$ and $\{C\}$ d -separate $\{A\}$ and $\{D\}$;

- $\mathbb{I}(A; B | C) \leq \mathbb{I}(A; B | D)$; In fact,

$$\begin{aligned} \mathbb{I}(A; B | C) &= \mathbb{I}(A; B | C, D) = \mathbb{I}(A; B, C | D) - \mathbb{I}(A; C | D) \\ &= \mathbb{I}(A; B | D) + \mathbb{I}(A; C | B, D) - \mathbb{I}(A; C | D) \\ &= \mathbb{I}(A; B | D) - \mathbb{I}(A; C | D) \leq \mathbb{I}(A; B | D), \end{aligned}$$

due to $\mathbb{I}(A; C | B, D) = 0$, because of

$$\begin{aligned} 0 &\leq \mathbb{I}(A; C | B, D) = \mathbb{I}(A; C, D | B) - \mathbb{I}(A; D | B) \\ &= \mathbb{I}(A; C | B) + \mathbb{I}(A; D | B, C) - \mathbb{I}(A; D | B) = -\mathbb{I}(A; D | B) \leq 0, \end{aligned}$$

since $\mathbb{I}(A; C | B) = 0$ (see Equation 17) and $\{B, C\}$ d -separates $\{A\}$ and $\{D\}$;

- $\mathbb{I}(A; B | C) \leq \mathbb{I}(A; B)$. In fact, by (17), we have $\mathbb{I}(A; C | B) = 0$. Thus,

$$\begin{aligned} \mathbb{I}(A; B) &= \mathbb{I}(A; B, C) - \mathbb{I}(A; C | B) = \mathbb{I}(A; C) + \mathbb{I}(A; B | C) - \mathbb{I}(A; C | B) \\ &= \mathbb{I}(A; C) + \mathbb{I}(A; B | C) \geq \mathbb{I}(A; B | C). \end{aligned}$$

(v) $\mathbb{I}(T; S | N) = \max_{S' \subseteq V \setminus N \setminus T} \mathbb{I}(T; S' | N)$: the proof is omitted.

(vi) $N \cup S$ is an Mb of T_1 in $V \setminus \{T_2\}$: By (17), we have

- $\mathbb{I}(A; C, D) \geq \mathbb{I}(A; C) > 0$;
- $\mathbb{I}(A; D | C) = 0$, since $\{C\}$ d -separates $\{A\}$ and $\{D\}$.

(vii) $N \cup S$ is an Mb of T_2 in $V \setminus \{T_1\}$: By (17), we have

- $\mathbb{I}(B; C, D) \geq \mathbb{I}(B; C) > 0$;
- $\mathbb{I}(B; D | C) = 0$, since $\{C\}$ d -separates $\{B\}$ and $\{D\}$.

(viii) S is an MBS; S_0 is an MbS (but not an MBS): $\mathbb{I}(A; B | C, D) = \mathbb{I}(A; B | C)$. In fact,

$$\mathbb{I}(A; B | C, D) = \mathbb{I}(A; B, D | C) - \mathbb{I}(A; D | C) = \mathbb{I}(A; B | C) + \mathbb{I}(A; D | B, C) - 0 = \mathbb{I}(A; B | C),$$

since both $\{B, C\}$ and $\{C\}$ d -separate $\{A\}$ and $\{D\}$. ■

Theorem 5 Let M_j be an MB of $T_j^{(d)}$ in $V \setminus T$ for $j = 1, \dots, t$. Then, $M \triangleq \cup_{j=1}^k M_j$ is an Mb of T . Further, M is an MB of T iff for any $X \in M$ there is some j such that $T_j^{(d)} \perp\!\!\!\perp X \mid M \setminus \{X\}$.

Proof Recall that T is the merged version of T , while $T^{(d)}$ is the dummy version of T ; all of them have the same MBs.

First, we have $T_j^{(d)} \perp\!\!\!\perp (V \setminus T) \setminus M_j \setminus \{T_j^{(d)}\} \mid M_j$ for $j = 1, \dots, t$. Considering $T_j^{(d)} \notin V \setminus T$, it follows that $T_j^{(d)} \perp\!\!\!\perp V \setminus T \setminus M_j \mid M_j$, which combined with the weak union property gives

$$T_j^{(d)} \perp\!\!\!\perp V \setminus T \setminus M \mid M, \quad j = 1, \dots, t, \quad (18)$$

since $M_j \subseteq M$. Putting $U \triangleq V \setminus T \setminus M$, the above independence statements imply

$$\mathbb{P}(T_j^{(d)} = 1, U = u \mid M) = \mathbb{P}(T_j^{(d)} = 1 \mid M) \mathbb{P}(U = u \mid M), \quad j = 1, \dots, t,$$

or equivalently, $\mathbb{P}(T = j, U = u \mid M) = \mathbb{P}(T = j \mid M) \mathbb{P}(U = u \mid M)$, meaning $T \perp\!\!\!\perp V \setminus T \setminus M \mid M$, and thus $T \perp\!\!\!\perp V \setminus T \setminus M \mid M$. This shows M is an Mb of T .

In what follows, we prove M is an MB of T if and only if, for any $X \in M$, $T_j^{(d)} \perp\!\!\!\perp X \mid M \setminus \{X\}$ holds for some j :

“ \Rightarrow ” Assume M is an MB of T . Suppose there is some variable X such that $T_j^{(d)} \perp\!\!\!\perp X \mid M \setminus \{X\}$ holds for any j . Then, by (18) and the contraction property, we get

$$T_j^{(d)} \perp\!\!\!\perp (V \setminus T) \setminus (M \setminus \{X\}) \mid M \setminus \{X\}, \quad j = 1, \dots, t. \quad (19)$$

Similar to the proof of the first conclusion, it can be readily proven that (19) implies

$$T \perp\!\!\!\perp (V \setminus T) \setminus (M \setminus \{X\}) \mid M \setminus \{X\},$$

meaning that M has a proper subset, $M \setminus \{X\}$, which is an Mb of T . This contradicts the minimality of M , and thus proves the necessity.

“ \Leftarrow ” Suppose M is not an MB of T (or T). Then, there is some $X \in M$ such that $T \perp\!\!\!\perp X \mid M \setminus \{X\}$. It follows that $\mathbb{P}(T = j, X = x \mid M \setminus \{X\}) = \mathbb{P}(T = j \mid M \setminus \{X\}) \mathbb{P}(X = x \mid M \setminus \{X\})$ holds for any $j = 1, \dots, t$. Or equivalently, we have

$$\mathbb{P}(T_j^{(d)} = 1, X = x \mid M \setminus \{X\}) = \mathbb{P}(T_j^{(d)} = 1 \mid M \setminus \{X\}) \mathbb{P}(X = x \mid M \setminus \{X\}). \quad (20)$$

Further, (20) indicates

$$\begin{aligned} \mathbb{P}(T_j^{(d)} = 0, X = x \mid M \setminus \{X\}) &= \mathbb{P}(X = x \mid M \setminus \{X\}) - \mathbb{P}(T_j^{(d)} = 1, X = x \mid M \setminus \{X\}) \\ &= [1 - \mathbb{P}(T_j^{(d)} = 1 \mid M \setminus \{X\})] \mathbb{P}(X = x \mid M \setminus \{X\}) \\ &= \mathbb{P}(T_j^{(d)} = 0 \mid M \setminus \{X\}) \mathbb{P}(X = x \mid M \setminus \{X\}). \end{aligned} \quad (21)$$

By (20) and (21), we get $T_j^{(d)} \perp\!\!\!\perp X \mid M \setminus \{X\}$, which contradicts $T_j^{(d)} \not\perp\!\!\!\perp X \mid M \setminus \{X\}$. This proves the sufficiency.

The proof is completed. ■

Theorem 6 (Correctness of IAMBS and KIAMBS) *Assume that T_2 satisfies the local composition property, and that all CI tests are correct. Then (i) IAMBS outputs an MB of $T_1 \cup T_2$; (ii) KIAMBS outputs an MB of $T_1 \cup T_2$ for any $K \in [0, 1)$.*

Proof Clearly, $N \cup S$ is an Mb of T_2 in $V \setminus T_1$ at the end of the growing phase of either IAMBS or KIAMBS under the local composition assumption, as in IAMB and KIAMB. Therefore, S is an MbS at the end of this stage. According to the last conclusion of Theorem 4, S is an MBS after it is refined. Finally, as a direct consequence of Lemma 4 (shown below), $N \cup S$ is an MB at the end of the algorithm, considering the process of refining N is similar to that of refining S . ■

Remark 2 *The following two statements hold: (a) violating local intersection implies violating adjacency faithfulness; (b) under the orientation faithfulness condition, violating local composition at the end of the first phase of IAMB or KIAMB or IAMBS or KIAMBS means violating adjacency faithfulness.*

Proof By Lemma 1, the violation of the local intersection property means information equivalence occurs; further, Lemeire et al. (2012) showed that information equivalence is one of the cases of violating adjacency faithfulness. Hence, the violation of local intersection is one of the violations of adjacency faithfulness.

Now, we show that the violation of local composition, which is present at the end of the first phase of IAMB or KIAMB, is also one of the violations of adjacency faithfulness under the *orientation faithfulness condition*.

In fact, let M be the output of the first phase of IAMB or KIAMB, but not an Mb of T . Without loss of generality, we assume $|T| = 1$ and $T = \{T\}$. Then, $T \perp\!\!\!\perp X \mid M$ holds for any $X \in V \setminus M \setminus \{T\}$ but $T \not\perp\!\!\!\perp V \setminus M \setminus \{T\} \mid M$. Considering that the set M_T composed of the parents, children, and spouses of T is an Mb of T , we have $M \not\supseteq M_T$. Thus, there is some $X \in M_T$ such that $X \notin M$. If X is a spouse of T , then all the children of T and X are not in M (if not so, $T \not\perp\!\!\!\perp X \mid M$ holds immediately following from the orientation faithfulness condition, and thus contradicts $T \perp\!\!\!\perp X \mid M$ since $X \notin M$). In this sense, we conclude that there is some node X adjacent to T such that $T \perp\!\!\!\perp X \mid M$. This means the adjacency faithfulness condition is violated.

In short words, both the violation of the local composition property (present at the end of the first phase of IAMB or KIAMB) and the violation of the local intersection property are the violations of adjacency faithfulness, under the orientation faithfulness condition. ■

Lemma 3 (a) *If there is $P \subseteq M_1 \setminus T_2$ such that $T_1 \perp\!\!\!\perp P \mid (N \setminus P) \cup T_2$, then $(N \setminus P) \cup T_2$ is an Mb of T_1 ; (b) If there is $Q \subseteq N \setminus P$ such that $T_1 \perp\!\!\!\perp Q \mid (N \setminus P \setminus Q) \cup T_2$ and $T_2 \perp\!\!\!\perp Q \mid (N \setminus P \setminus Q) \cup T_1$, then $(N \setminus P \setminus Q) \cup T_2$ is an Mb of T_1 , and $(N \setminus P \setminus Q) \cup T_1$ is an Mb of T_2 .*

Proof Considering that $N \cup T_2$ is an Mb of T_1 , we have $T_1 \perp\!\!\!\perp V \setminus (N \cup T_2) \setminus T_1 \mid (N \setminus P) \cup T_2 \cup P$, which combined with $T_1 \perp\!\!\!\perp P \mid (N \setminus P) \cup T_2$ implies $T_1 \perp\!\!\!\perp V \setminus [(N \setminus P) \cup T_2] \setminus T_1 \mid (N \setminus P) \cup T_2$, in view of the contraction property. The first conclusion is proved.

For convenience, we denote now $N_1 \triangleq M_1 \setminus T_2$. To show the second conclusion, we note that $P \subseteq N_1$, so $(N \setminus P) \cup T_1$ is an Mb of T_2 . It follows that: (i) $T_1 \perp\!\!\!\perp V \setminus [(N \setminus P) \cup T_2] \setminus T_1 \mid (N \setminus P) \cup T_2$, which combined with $T_1 \perp\!\!\!\perp Q \mid (N \setminus P \setminus Q) \cup T_2$ gives $T_1 \perp\!\!\!\perp V \setminus [(N \setminus P \setminus Q) \cup T_2] \setminus T_1 \mid (N \setminus P \setminus Q) \cup T_2$; and (ii) $T_2 \perp\!\!\!\perp V \setminus [(N \setminus P) \cup T_1] \setminus T_2 \mid (N \setminus P) \cup T_1$, which combined with $T_2 \perp\!\!\!\perp Q \mid (N \setminus P \setminus Q) \cup T_1$ yields $T_2 \perp\!\!\!\perp V \setminus [(N \setminus P \setminus Q) \cup T_1] \setminus T_2 \mid (N \setminus P \setminus Q) \cup T_1$. The second conclusion is also proved. ■

Lemma 4 Let T_i be a subset of V with an Mb M_i for $i = 1, 2$, and $S \subseteq V \setminus N \setminus T$ be an MBS of T to N , with $T = T_1 \cup T_2$ and $N = (M_1 \cup M_2) \setminus T$. Assume N_0 be a subset of N such that $N_0 \cup S$ is an Mb of T . Then $N_0 \cup S$ is an MB of T if and only if $T \perp\!\!\!\perp Y | (N_0 \setminus \{Y\}) \cup S$ holds for any $Y \in N_0$.

Proof By the definition of MBS, $N \cup R$ and thus $N_0 \cup R$ will never be an Mb of T for any $N_0 \subseteq N$ and $R \subsetneq S$, in view of the weak union property.

- Necessity: Suppose there is some $Y \in N_0$ such that $T \perp\!\!\!\perp Y | (N_0 \setminus \{Y\}) \cup S$. By the precondition that $N_0 \cup S$ is an Mb of T , we have $T \perp\!\!\!\perp V \setminus (N_0 \cup S) \setminus T | [(N_0 \setminus \{Y\}) \cup S] \cup \{Y\}$. These two relationships combined with the contraction property imply

$$T \perp\!\!\!\perp V \setminus [(N_0 \setminus \{Y\}) \cup S] \setminus T | (N_0 \setminus \{Y\}) \cup S,$$

or equivalently, $(N_0 \setminus \{Y\}) \cup S$ is an Mb of T . This contradicts that $N_0 \cup S$ is an MB of T .

- Sufficiency: Suppose $N_0 \cup S$ is not an MB of T , that is, there is some $N'_0 \subsetneq N_0$ such that $N'_0 \cup S$ is an Mb of T . Take any given variable, Y , in $N_0 \setminus N'_0$. It can be shown that $T \perp\!\!\!\perp Y | (N_0 \setminus \{Y\}) \cup S$, which leads to a contradiction. Hence, $N_0 \cup S$ is an MB of T .

The proof is completed. ■

Appendix C. Improving the Log-Likelihood Ratio Test

In Subsection 4.3, we mentioned that the X^2 or G^2 test is suitable only for cases of small $|T|$, and then summarized some improving methods proposed in the literature (Lawley, 1956; Hosmane, 1986, 1987, 1990; Brin et al., 1997; Silverstein et al., 1998; Aliferis et al., 2010b). However, we need more suitable CI testing methods when working on the MB discovery problem for multiple targets. In this appendix, we discuss a practical way of improving the G^2 test by damping the number of degrees of freedom for the G^2 statistic.

Consider the G^2 statistic, $G^2(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \triangleq 2n \cdot \mathbb{I}_D(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$, which approximates to the chi-square variate with $r \triangleq (r_X - 1)(r_Y - 1)r_Z$ degrees of freedom, namely $\chi^2(r)$, where r_ξ represents the number of configurations for ξ (de Campos, 2006, p. 2158).

Theoretically, $G^2(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$ is a reasonable statistic for testing the hypothesis “ $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}$ ” when n is large enough. Unfortunately, this precondition is practically hard to be valid in many situations (Cochran, 1954; Yaramakala, 2004; Bromberg and Margaritis, 2009) due to the following reason: Let $\mathbf{X} = \{X_{i_1}, \dots, X_{i_x}\}$, $\mathbf{Y} = \{X_{j_1}, \dots, X_{j_y}\}$, and $\mathbf{Z} = \{X_{k_1}, \dots, X_{k_z}\}$, in which each variable X_ℓ takes r_ℓ values. Then $r = (\prod_{\ell=1}^x r_{i_\ell} - 1)(\prod_{\ell=1}^y r_{j_\ell} - 1)(\prod_{\ell=1}^z r_{k_\ell})$, which is exponential with respect to x , y , and z . On the one hand, by the Wilson-Hilferty approximation for $\chi^2_\alpha(r)$ (de Campos, 2006; Gao, 2005), we obtain $\chi^2_\alpha(r) \approx c_{\alpha,r} r$, in which $c_{\alpha,r} \triangleq (1 - 2/(9r) + \sqrt{2/(9r)} z_\alpha)^3$ is a bit larger than 1, with z_α being the upper α -quantile of the standard normal distribution; on the other hand, we can show $\mathbb{I}_D(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \leq \log_2 r_{X,Y}$ with $r_{X,Y} \triangleq \min\{\prod_{\ell=1}^x r_{i_\ell}, \prod_{\ell=1}^y r_{j_\ell}\}$. It follows that

$$p(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = \mathbb{P}\{\chi^2(r) \geq 2n \cdot \mathbb{I}_D(\mathbf{X}; \mathbf{Y} | \mathbf{Z})\} \geq \mathbb{P}\{\chi^2(r) \geq 2n \cdot \log_2 r_{X,Y}\}.$$

Suppose we are doing a G^2 test for the *false hypothesis* “ $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}$ ” (i.e., the truth is $\mathbf{X} \not\perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}$). Then, at least $\frac{\chi^2_\alpha(r)}{2 \log_2 r_{X,Y}} \approx \frac{c_{\alpha,r} r}{2 \log_2 r_{X,Y}} = O\left(\frac{r}{\log_2 r_{X,Y}}\right)$ instances are required if we expect the statistical

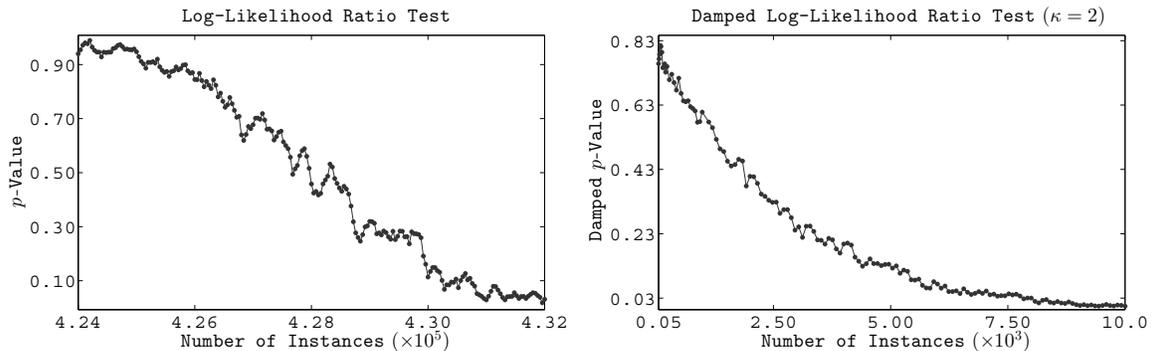


Figure 20: p -value/damped p -value versus the number of instances, n : the left subfigure illustrates why a very large n may still not be “large enough” for making a correct decision about the *false hypothesis* “ $X_1 \perp\!\!\!\perp Y_1 | Z_1$ ” based on the G^2 test: at least 4.31×10^5 instances are required; while the right illustrates why the damped G^2 test is suitable for testing the same *false hypothesis*: about 8000 instances are sufficient.

decision can be correctly made with the significance level α (or equivalently, $p(X; Y | Z) \leq \alpha$). In many practical situations, however, n may be far smaller than the required number of sample instances with the magnitude of at least $O(r / \log_2 r_{X,Y})$, recalling that r is exponential with respect to x , y , and z , especially when too many variables are involved. In this case, the statistical decision made for the hypothesis will be wrong.

Taking the ALARM network presented in Figure 4 for example, we put

$$X_1 = \{X_{36}\}, Y_1 = \{X_{11}, X_{34}, X_{35}, X_{37}\}, \text{ and } Z_1 = \{X_4, X_{14}, X_{15}, X_{16}, X_{18}, X_{21}, X_{22}, X_{31}\};$$

then we compute the p -value versus the number of instances from 1000 to 1,000,000. The results are drawn in Figure 20 (averaged over 10 different samples with the same size). Note that the truth is $X_1 \not\perp\!\!\!\perp Y_1 | Z_1$ since Y_1 is an Mb of X_1 . By the figure, the CI test for the *false hypothesis* “ $X_1 \perp\!\!\!\perp Y_1 | Z_1$ ” is not correct unless at least $n_{\min} \approx 4.31 \times 10^5$ sample instances are available. It is mentioned that, in this example, $r / \log_2 r_{X_1, Y_1} \approx 2.42 \times 10^5$.

In short words, the precondition, “when n is large enough”, for the theoretical assertion that “ G^2 is a reasonable statistic for CI testing” may be hard to be guaranteed in practice because the above analysis and the numerical example indicate that a seemingly very large n may still not be “large enough”. The problem is then how to improve on the G^2 test.

Observe that, for the G^2 test, the major reason for failing to make a correct statistical decision on CI testing is that the theoretical value of r is far larger than its data-driven value, due to the null cells frequently existing in the multi-contingency tables of X and Y given Z (e.g., Yaramakala, 2004, p. 34). In other words, the linear increase of n is hard to exponentially bring null cells into valid cells. Hence, a feasible way of improving the log-likelihood ratio G^2 test is to damp the increase of r such that the unmatched behaviours of n and r can get alleviated to a certain degree. Mathematically, we replace the theoretical value of r in $p(X; Y | Z)$ with its a damped version, $g_{n,\kappa}(r)$, defined in (6), where $\kappa > 0$ is a constant, based on which $\frac{n}{\kappa}$ measures the amount of valid cells that n sample

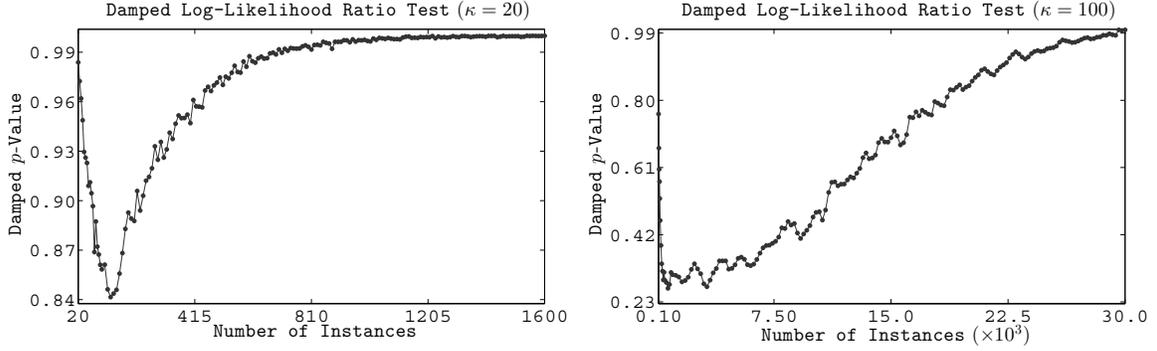


Figure 21: Damped p -value versus the number of instances: κ is taken as 20 and 100, respectively.

instances can support. It is easily seen that $g_{n,\kappa}(\cdot)$ possesses the following properties, which interpret the reasonability of employing such a damping procedure in the G^2 test:

- $g_{n,\kappa}(r)$ is monotonically increasing versus n for given r , and $\lim_{n \rightarrow +\infty} g_{n,\kappa}(r) = r$. This means more instances may generate more valid cells in the multi-contingency tables of \mathbf{X} and \mathbf{Y} given \mathbf{Z} , and all the theoretical degrees of freedom are valid when n is large enough.
- $g_{n,\kappa}(r)$ is monotonically increasing versus r for given n , and $\lim_{r \rightarrow +\infty} g_{n,\kappa}(r) = \frac{n}{\kappa}$. This means a larger r should correspond to a larger $g_{n,\kappa}(r)$, but not exceeding the supporting capacity of the data.
- For sufficient data, the damping function $g_{n,\kappa}(\cdot)$ only plays a little role; while for insufficient data, it trades off the theoretical r and the supporting capacity of the data.

For convenience, we call the resulted p -value, denoted by $p_g(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$ instead of $p(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$, and the resulted testing method to be the *damped p -value* and the *damped log-likelihood ratio test* (or damped G^2 test). Further, we use the negative damped p -value, $f_D^{(3)}$ defined in (7), as the association function. It is mentioned here that the damped G^2 test approximately degenerates into the ordinary G^2 test when taking κ as a very small positive number.

The damped G^2 test may be more suitable than the ordinary G^2 test when too many variables are involved in the conditional set. We implement this testing method (by taking κ as 2, 3, \dots , 10, respectively) on the *false hypothesis* “ $\mathbf{X}_1 \perp\!\!\!\perp \mathbf{Y}_1 | \mathbf{Z}_1$ ”, and find that the correct decision “ $\mathbf{X}_1 \not\perp\!\!\!\perp \mathbf{Y}_1 | \mathbf{Z}_1$ ” is always made for $\kappa \geq 3$, even when the number of instances is smaller than 1000. For the case of $\kappa = 2$, we present the results in the right subfigure of Figure 20, from which it is seen that, for the damped G^2 test, about 8000 sample instances are sufficient to make the correct decision. Note also that the ordinary G^2 test needs at least 4.31×10^5 instances.

However, the damped G^2 test may also face a potential danger: it may excessively damp the theoretical value of r if a too large κ is inappropriately used. Here, “excessively damping r ” means that a too large value of κ will lead to a too small $g_{n,\kappa}(r)$ such that the damped G^2 test incorrectly reject a *true hypothesis* “ $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}$ ”. To illustrate this explanation, we put

$$\mathbf{X}_2 = \{X_{21}\}, \quad \mathbf{Y}_2 = \{X_{11}, X_{34}, X_{35}, X_{36}, X_{37}\}, \quad \text{and} \quad \mathbf{Z}_2 = \{X_{15}, X_{19}, X_{20}, X_{22}, X_{29}\}$$

from the ALARM network. The truth is $\mathbf{X}_2 \perp\!\!\!\perp \mathbf{Y}_2 | \mathbf{Z}_2$ since \mathbf{Z}_2 is an Mb of \mathbf{X}_2 . Now, use damped G^2 to test the *true hypothesis* “ $\mathbf{X}_2 \perp\!\!\!\perp \mathbf{Y}_2 | \mathbf{Z}_2$ ” by taking κ as 3, 4, \dots , 10, 20, 100, respectively. All

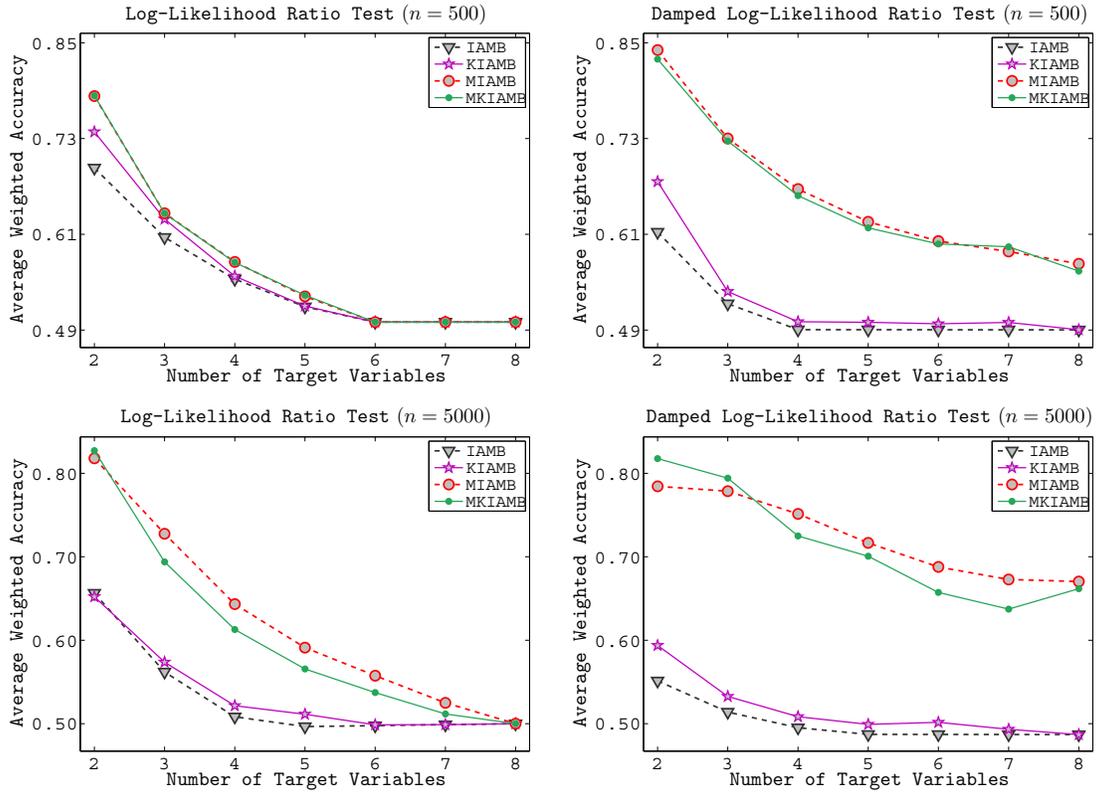


Figure 22: An illustration on why the damped G^2 test is more suitable than the ordinary G^2 test for the problem that involves too many variables, by virtue of the Pigs BN: the G^2 test no longer works when there are too many variables are involved, while the damped G^2 test remains valid in all considered cases.

the decisions are correctly made. However, a too large κ may be more apt to yield a relatively small damped p -value although it is still larger than α , as shown by Figure 21: 0.84 for the case of $\kappa = 20$ and only 0.27 for the case of $\kappa = 100$. Hence, to avoid the potential danger of excessively damping r , we conservatively recommend to take κ from the interval $[3, 10]$ in practice. In our benchmarking study, we employ $\kappa = 5$, which is large enough for testing *false hypotheses* and small enough for testing *true hypotheses*.

To further illustrate why the damped G^2 test is more suitable than the ordinary G^2 when working on a problem that involves too many variables, we make experiments on the six synthetic BNs based on IAMB, KIAMB, MIAMB, and MKIAMB.

For each algorithm, the G^2 test and the damped G^2 test are implemented for CI testing. Accordingly, the association functions, $f_D^{(1)}$ and $f_D^{(3)}$ defined by (5) and (7) are used. Figure 22 presents the results of the Pigs network, in which the left two are based on the ordinary G^2 test while the right two are based on the damped G^2 test. As seen, for the case of $n = 500$ the G^2 test becomes invalid when the target, \mathbf{T} , contains 6 or more variables, and for the case of $n = 5000$ this method no longer works when $|\mathbf{T}| = 8$. In comparison, the damped G^2 test is suitable for all cases. The results of the other five BNs show similar conclusions.

Appendix D. Used Acronyms

BFMB	breadth first search of Markov boundary algorithm (Fu and Desmarais, 2007).
BN	Bayesian network.
CI	conditional independence.
CMI	conditional mutual information.
CT-support	contingency table support: a set of items S has CT-support s at the $t\%$ level if at least $t\%$ of the cells in the contingency table for S have value s (Silverstein et al., 1998).
DAG	directed acyclic graph.
E-partition	equivalent partition (Lemeire, 2007): a relation $\mathfrak{R} \subset X \otimes Y$ defines an E-partition in Y_{dom} to a partition of X_{dom} , if: (i) $\neg(x_2 \mathfrak{R} y_1)$ holds for any $x_1, x_2 \in X_{\text{dom}}$ belonging to different partitions and for any $y_1 \in Y_{\text{dom}}$ with $x_1 \mathfrak{R} y_1$; and (ii) for every $X_{\text{dom}}^{(k)}$, there exist $x_1 \in X_{\text{dom}}^{(k)}$ and $y_1 \in Y_{\text{dom}}$ such that $x_1 \mathfrak{R} y_1$.
FS	feature selection.
GLL	generalized local learning: an algorithmic framework for local causal discovery and FS proposed by Aliferis et al. (2010a).
GS	grow-shrink algorithm (Margaritis and Thrun, 1999, 2000).
HITON	an MB discovery algorithm, pronounced hee-tón, from the Greek <i>Χιτώνας</i> , for “cover”, “cloak”, or “blanket” (Aliferis et al., 2003).
IAMB	incremental association Markov boundary algorithm (Tsamardinos et al., 2003); see Algorithm 3 for details.
IAMBS	an IAMB-based Markov boundary supplementary algorithm, outputting an MB for multiple targets (Algorithm 1).
KIAMB	a stochastic variant of IAMB (Peña et al., 2007); see Algorithm 3 for details.
KIAMBS	an KIAMB-based Markov boundary supplementary algorithm, outputting an MB for multiple targets (Algorithm 1).
KS	Koller-Sahami algorithm (Koller and Sahami, 1996).
LibSVM	a library for support vector machines contributed by Chang and Lin (2011).
Mb	Markov blanket: we call M an Mb of T if $T \perp\!\!\!\perp V \setminus M \setminus T \mid M$ (Definition 1).
MB	Markov boundary: an MB of T is any Mb such that none of its proper subsets is an Mb of T (Definition 1).
MbS	Markov blanket supplementary: we call S an MbS of T to N , if $N \cup S$ is an Mb of T (Definition 3).
MBS	Markov boundary supplementary: an MBS is any MbS such that none of its proper subsets is an MbS (Definition 3).
MIAMB	an IAMB and IAMBS-based algorithm, outputting an MB for multiple targets (see Algorithm 2 for details).
MKIAMB	an KIAMB and KIAMBS-based algorithm, outputting an MB for multiple targets (Algorithm 2).
MMMB	max-min Markov boundary algorithm (Tsamardinos et al., 2006).

PCMB	parents and children based Markov boundary algorithm (Peña et al., 2007).
RF	random forest algorithm.
RT	running time: the single CPU time implemented on an Intel i7-3612QM 2.1 GHz and Windows 7 with 64 bits.
SVM	support vector machine (in one-against-one approach).
T-partition	target partition (Lemeire, 2007): the domain, X_{dom} , of X can be partitioned into some disjoint subsets $X_{\text{dom}}^{(k)}$ for which $\mathbb{P}(T \mathbf{x})$ is the same for all $\mathbf{x} \in X_{\text{dom}}^{(k)}$. This is called the T-partition of X_{dom} with respect to T .
WA	weighted accuracy: WA is the average of the rate of true members and that of true nonmembers of an MB with respect to the truth.
WP	weighted precision: WP is the average of the rate of true members and that of true nonmembers of an MB with respect to the test.

References

- Constantin F Aliferis, Ioannis Tsamardinos, and Alexander Statnikov. Hiton: a novel Markov blanket algorithm for optimal variable selection. In *AMIA 2003 Annual Symposium Proceedings*, pages 21–25. American Medical Informatics Association, 2003.
- Constantin F Aliferis, Alexander Statnikov, Ioannis Tsamardinos, Subramani Mani, and Xenofon D Koutsoukos. Local causal and Markov blanket induction for causal discovery and feature selection for classification part I: Algorithms and empirical evaluation. *Journal of Machine Learning Research*, 11:171–234, 2010a.
- Constantin F Aliferis, Alexander Statnikov, Ioannis Tsamardinos, Subramani Mani, and Xenofon D Koutsoukos. Local causal and Markov blanket induction for causal discovery and feature selection for classification part II: Analysis and extensions. *Journal of Machine Learning Research*, 11:235–284, 2010b.
- Ingo A Beinlich, H J Suermondt, R Martin Chavez, and Gregory F Cooper. The ALARM monitoring system: A case study with two probabilistic inference techniques for belief networks. In *Second European Conference on Artificial Intelligence in Medicine*, pages 247–256, London, 1989. Springer-Verlag.
- Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL <https://doi.org/10.1023/A:1010933404324>.
- Sergey Brin, Rajeev Motwani, and Craig Silverstein. Beyond market baskets: Generalizing association rules to correlations. In *SIGMOD'97 Proceedings of the 1997 ACM SIGMOD International Conference on Management of Data*, volume 26, pages 265–276. ACM, 1997.
- Facundo Bromberg and Dimitris Margaritis. Improving the reliability of causal discovery from small data sets using argumentation. *Journal of Machine Learning Research*, 10:301–340, 2009.
- Gavin Brown, Adam Pocock, Ming-Jie Zhao, and Mikel Luján. Conditional likelihood maximisation: a unifying framework for information theoretic feature selection. *Journal of Machine Learning Research*, 13:27–66, 2012. (Version: MIToolbox-2.0).

- Chih-Chung Chang and Chih-Jen Lin. LibSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Jie Cheng, Russell Greiner, Jonathan Kelly, David Bell, and Weiru Liu. Learning Bayesian networks from data: An information-theory based approach. *Artificial Intelligence*, 137(1–2):43–90, 2002.
- William G Cochran. Some methods for strengthening the common χ^2 tests. *Biometrics*, 10(4):417–451, 1954.
- Thomas M Cover and Joy A Thomas. *Elements of Information Theory (Second Edition)*. John Wiley and Sons, 2006.
- Noel Cressie and Timothy RC Read. Pearson’s X^2 and the loglikelihood ratio statistic G^2 : a comparative review. *International Statistical Review*, 57(1):19–43, 1989.
- Rónán Daly, Qiang Shen, and Stuart Aitken. Learning Bayesian networks: Approaches and issues. *The Knowledge Engineering Review*, 26(2):99–157, 2011.
- Luis M de Campos. A scoring function for learning Bayesian networks based on mutual information and conditional independence tests. *Journal of Machine Learning Research*, 7(Oct):2149–2187, 2006.
- Shunkai Fu and Michel Desmarais. Local learning algorithm for Markov blanket discovery. In *AI 2007: Advances in Artificial Intelligence*, pages 68–79. Springer Berlin Heidelberg, 2007.
- Shunkai Fu and Michel C Desmarais. Markov blanket based feature selection: a review of past decade. In *Proceedings of the World Congress on Engineering*, 2010.
- Huixuan Gao. *Statistics Computation*. Peking University Press, Beijing, 2005.
- Balakrishna Hosmane. Improved likelihood ratio test for multinomial goodness of fit. *Communications in Statistics-Theory and Methods*, 16(11):3185–3198, 1987.
- Balakrishna S Hosmane. Smoothing of likelihood ratio statistic for equiprobable multinomial goodness-of-fit. *Annals of the Institute of Statistical Mathematics*, 42(1):133–147, 1990.
- BS Hosmane. Improved likelihood ratio tests and Pearson chi-square tests for independence in two dimensional contingency tables. *Communications in Statistics-Theory and Methods*, 15(6):1875–1888, 1986.
- James Kennedy and Russell C Eberhart. Particle swarm optimization. In *Proceedings of 1995 IEEE International Conference on Neural Networks*, volume 4, pages 1942–1948, Perth, 1995.
- James Kennedy and Russell C Eberhart. A discrete binary version of the particle swarm algorithm. In *IEEE International Conference on Systems, Man, and Cybernetics*, volume 5, pages 4104–4108, Orlando, 1997.
- Daphne Koller and Mehran Sahami. Toward optimal feature selection. In *Thirteen International Conference in Machine Learning*. Stanford InfoLab, 1996.

- D.N. Lawley. A general method for approximating to the distribution of likelihood ratio criteria. *Biometrika*, 43(3–4):295–303, 1956.
- Jan Lemeire. *Learning Causal Models of Multivariate Systems and the Value of it for the Performance Modeling of Computer Programs*. ASP/VUBPRESS/UPA, PhD thesis, 2007.
- Jan Lemeire, Stijn Meganck, Francesco Cartella, and Tingting Liu. Conservative independence-based causal structure learning in absence of adjacency faithfulness. *International Journal of Approximate Reasoning*, 53(9):1305–1325, 2012.
- Dimitris Margaritis and Sebastian Thrun. Bayesian network induction via local neighborhoods. Technical Report CMU-CS-99-134, 1999.
- Dimitris Margaritis and Sebastian Thrun. Bayesian network induction via local neighborhoods. In *Advances in Neural Information Processing Systems*, volume 12, pages 505–511. Morgan Kaufmann, 2000.
- Kevin P. Murphy. *Bayes Net Toolbox for Matlab*, 2007. (Version: FullBNT-1.0.7).
- Richard E Neapolitan. *Learning Bayesian Networks*. Upper Saddle River: Prentice Hall, 2004.
- Pekka Parviainen and Mikko Koivisto. Finding optimal Bayesian networks using precedence constraints. *Journal of Machine Learning Research*, 14(1):1387–1415, 2013.
- Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco: Morgan Kaufmann, 1988.
- Jean-Philippe Pellet and André Elisseeff. Using Markov blankets for causal structure learning. *Journal of Machine Learning Research*, 9:1295–1342, 2008.
- Jose M Peña, Roland Nilsson, Johan Björkegren, and Jesper Tegnér. Towards scalable and data efficient learning of Markov boundaries. *International Journal of Approximate Reasoning*, 45(2): 211–232, 2007.
- Joseph Ramsey, Jiji Zhang, and Peter L Spirtes. Adjacency-faithfulness and conservative causal inference. In *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence (UAI-2006)*, pages 401–408, 2006.
- Johannes Rauh, Nils Bertschinger, Eckehard Olbrich, and Jurgen Jost. Reconsidering unique information: Towards a multivariate information decomposition. In *2014 IEEE International Symposium on Information Theory (ISIT)*, pages 2232–2236. IEEE, 2014.
- Federico Schlüter. A survey on independence-based Markov networks learning. *Artificial Intelligence Review*, 42:1069–1093, 2014.
- Marco Scutari. Learning Bayesian networks with the bnlearn R package. *Journal of Statistical Software*, 35:1–22, 2010.
- Craig Silverstein, Sergey Brin, and Rajeev Motwani. Beyond market baskets: Generalizing association rules to dependence rules. *Data Mining and Knowledge Discovery*, 2(1):39–68, 1998.

- Alexander Statnikov and Constantin F Aliferis. Analysis and computational dissection of molecular signature multiplicity. *PLoS Computational Biology*, 6(5):e1000790, 2010.
- Alexander Statnikov, Nikita I. Lytkin, Jan Lemeire, and Constantin F Aliferis. Algorithms for discovery of multiple Markov boundaries. *Journal of Machine Learning Research*, 14(1):499–566, 2013.
- Ioannis Tsamardinos and Constantin F Aliferis. Towards principled feature selection: Relevancy, filters and wrappers. In *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*, 2003.
- Ioannis Tsamardinos, Constantin F Aliferis, and Alexander Statnikov. Algorithms for large scale Markov blanket discovery. In *Proceedings of the Sixteenth International Florida Artificial Intelligence Research Society Conference (FLAIRS)*, pages 376–381, 2003.
- Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing Bayesian network structure learning algorithm. *Machine learning*, 65(1):31–78, 2006.
- Sandeep Yaramakala. *Fast Markov Blanket Discovery*. MS thesis, 2004.
- Lianwen Zhang and Haipeng Guo. *Introduction to Bayesian Networks*. Science Press, Beijing, 2006.