Accurate Markov Boundary Discovery for Causal Feature Selection

Xingyu Wu, Bingbing Jiang[®], Kui Yu[®], Chunyan Miao[®], and Huanhuan Chen[®], Senior Member, IEEE

Abstract—Causal feature selection has achieved much atten-2 tion in recent years, which discovers a Markov boundary (MB) 3 of the class attribute. The MB of the class attribute implies local 4 causal relations between the class attribute and the features, 5 thus leading to more interpretable and robust prediction mod-6 els than the features selected by the traditional feature selection 7 algorithms. Many causal feature selection methods have been 8 proposed, and almost all of them employ conditional indepen-9 dence (CI) tests to identify MBs. However, many datasets from 10 real-world applications may suffer from incorrect CI tests due 11 to noise or small-sized samples, resulting in lower MB discovery 12 accuracy for these existing algorithms. To tackle this issue, in 13 this article, we first introduce a new concept of PCMasking to 14 explain a type of incorrect CI tests in the MB discovery, then 15 propose a cross-check and complement MB discovery (CCMB) 16 algorithm to repair this type of incorrect CI tests for accurate 17 MB discovery. To improve the efficiency of CCMB, we further 18 design a pipeline machine-based CCMB (PM-CCMB) algorithm. 19 Using benchmark Bayesian network datasets, the experiments 20 demonstrate that both CCMB and PM-CCMB achieve signif-21 icant improvements on the MB discovery accuracy compared 22 with the existing methods, and PM-CCMB further improves the 23 computational efficiency. The empirical study in the real-world 24 datasets validates the effectiveness of CCMB and PM-CCMB 25 against the state-of-the-art causal and traditional feature selection 26 algorithms.

Index Terms—Bayesian network (BN), causal feature selection,
 Markov boundary (MB), PCMasking.

I. INTRODUCTION

AUSAL feature selection is to identify a Markov boundary (MB) of a class attribute for building accurate prediction models. The MB was first defined and discussed by Pearl in a Bayesian network (BN) [1]. Under the faithfulness

Manuscript received May 5, 2019; revised August 17, 2019; accepted August 26, 2019. This work was supported in part by the National Key Research and Development Program of China under Grant 2016YFB1000905, in part by the Beijing Municipal Science and Technology Commission under Grant Z171100000117017, and in part by the National Natural Science Foundation of China under Grant 61876206, Grant 91846111, and Grant 91746209. This article was recommended by Associate Editor X. Wang. (Corresponding author: Huanhuan Chen.)

X. Wu, B. Jiang, and H. Chen are with the School of Computer Science and Technology, University of Science and Technology of China, Hefei 230027, China (e-mail: xingyuwu@mail.ustc.edu.cn; jiangbb@mail.ustc.edu.cn; hchen@ustc.edu.cn).

K. Yu is with the School of Computer Science and Information Engineering, Hefei University of Technology, Hefei 230601, China (e-mail: yukui@hfut.edu.cn).

C. Miao is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798 (e-mail: ascymiao@ntu.edu.sg). Color versions of one or more of the figures in this article are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TCYB.2019.2940509

assumption (refer to Definition 5 in Section III), the MB of a variable in a BN consists of its parents, children, and spouses (the other parents of the children of the variable), and given the MB of a variable, all other variables will be independent of this variable [1]. Thus, the MB provides a complete picture of the local causal structure around a variable [2]. In addition, in theory, the MB of the class attribute is the optimal solution of the feature selection problem [3], [4].

Recent years have witnessed the proliferation of the causal feature selection methods since they can select features not only predictive but also causal informative. Existing algorithms can be roughly divided into two different types. The first type is to directly discover the MB of a target variable of interest, which sacrifices the MB discovery accuracy to improve the computational efficiency. The early causal feature selection methods, such as the growth and shrink algorithm (GS) [5] and the Koller-Sahami (KS) [6] algorithms, belong to the first type. The later methods, incremental association MB (IAMB) and its variants [7], improve the GS by reordering the variables each time the MB set changes. However, IAMB and its variants require large amount of data to guarantee the accuracy. To solve this problem, the second type of algorithms is proposed by employing a divideand-conquer strategy to improve the MB discovery accuracy. Min-max MB (MMMB) [8] is the first divide-and-conquerbased method, later algorithms, such as HITON-MB [9] and parents-children-based MB (PCMB) [10], are improved on MMMB, which first find the parents and children (PC) of a target, and then identify the spouse (SP) of the target. MBOR [11] combines the two types of methods, which employs the first type of method to obtain an initial MB first and then finds more MB variables with the divide-and-conquer strategy.

To improve the MB discovery accuracy, the existing causal feature selection algorithms mainly focus on how to remove the false positives during the MB search process, but rarely consider the true positives discarded due to incorrect conditional independence (CI) tests, leading to low true positive discovery accuracy, especially in the presence of insufficient or noise data samples. For example in Fig. 1, by conducting the experiment on a benchmark Alarm BN with different sample sizes, we found that the recall of the existing causal feature selection algorithms is much smaller than their precision. Specifically, given 1000 samples, the average precision of these algorithms is 0.92, but the average recall is only 0.81 (a more detailed comparison on accuracy can be found in Fig. 5 of Section VI). The low recall value makes the existing causal

65

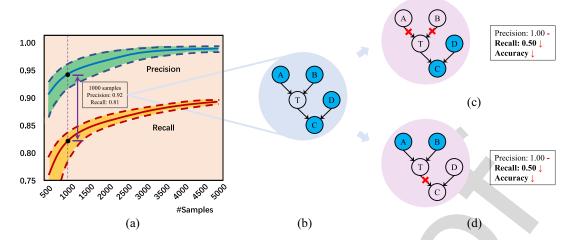


Fig. 1. Existing causal feature selection algorithms have higher precision but lower recall. A series of experiments was conducted on Alarm BN with different sample scale and averaging over several state-of-the-art algorithms (IAMB, PCMB, MBOR, and STMB). (a) Average precision and recall variation curves with respect to the number of samples. (b)-(d) Take a DAG as an example to illustrate the higher precision and lower recall of MB discovery when there exists the PCMasking phenomenon. (b) Correct result and highlight the true MB of T in blue, that is, parents A and B, child C, and spouse D. (c) and (d) The two wrong cases where the PCMasking phenomenon occurs. The wrong MBs of T in (c) and (d) are highlighted in blue, and the red "X" symbol denotes the independence relations between T and its parents or children.

81 feature selection algorithms ineffective for accurate prediction practical applications.

Motivated by the above issue, we experimented on different 83 84 benchmark BN datasets using existing causal feature selection 85 algorithms to explore why some true positives are discarded. We found a type of incorrect CI tests that makes PC of a tar-87 get mask each other, and we call it PCMasking. Specifically, 88 PCMasking denotes a target and its children may be indepen-89 dent conditioning on its parents and vice-versa. For example, e use a directed acyclic graph (DAG) in Fig. 1(b) to illustrate $_{91}$ the PCMasking phenomenon. Assuming that the target T and 92 its parents $\{A, B\}$ are independent conditioning on its child C $_{93}$ [as Fig. 1(c)] and, meanwhile, T and C are independent condi-94 tioning on $\{A, B\}$ [as Fig. 1(d)]. The incorrect tests will make $\{A, B\}$ or C not be added to the final output of the existing ₉₆ algorithms. In Fig. 1(b), the MB of T should be $\{A, B, C, D\}$. 97 However, if we apply existing algorithms to the dataset and 98 the type of incorrect CI tests occurs, the output of the algo-99 rithms is $\{A, B\}$ [as highlighted in Fig. 1(c)] or $\{C, D\}$ [as 100 highlighted in Fig. 1(d)] due to the PCMasking phenomenon [highlighted with red "x" in Fig. 1(c) and (d)].

Besides empirical analysis, we also analyze the mecha-103 nism of the PCMasking phenomenon from the perspective of information theory, and further find that PCMasking phe-105 nomenon breaks the symmetry between the PC variables, which leads to some true PC variables discarded by the existing algorithms. Therefore, if there is no extra strategy to deal with the PCMasking phenomenon, then the MB discovery methods will ignore some direct causes (parents) and direct 110 effects (children), resulting in the performance degradation. Furthermore, since the algorithms complete the PC discov-112 ery first and then search for the spouses based on the PC 113 set, incomplete PC set will cause cascading errors in the SP discovery. However, the PCMasking phenomenon has attracted 115 little attention, resulting in ineffectiveness of the existing 116 causal feature selection algorithms on real-world datasets. To tackle this issue, the main contributions of this article are 117 summarized as follows.

- 1) We formally present a new concept, called PCMasking, 119 to describe a type of incorrect CI tests in the MB dis- 120 covery process, and theoretically analyze the mechanism 121 behind this type of CI tests.
- 2) Based on the theoretical analyses, we propose the 123 cross-check and complement MB discovery (CCMB) 124 algorithm to tackle the PCMasking phenomenon and 125 improve the true-positive discovery accuracy. Moreover, 126 to improve the computational efficiency of CCMB, we 127 further design a PM-CCMB algorithm, which is more 128 accurate and efficient compared with all other MB 129 discovery algorithms.
- 3) We conduct a series of experiments on synthetic and 131 real-world datasets, to validate the effectiveness and 132 efficiency of the proposed algorithms against the state- 133 of-the-art causal and traditional feature selection algo- 134

The remainder of this article is organized as follows. 136 Section II reviews the related work and Section III intro- 137 duces the basic notations and definitions. In Section IV, we 138 propose the new concept of PCMasking and explain its mech- 139 anism. The proposed CCMB and PM-CCMB algorithms are 140 described in Sections V. The experimental results and analy- 141 ses are presented in Section VI. Finally, Section VII concludes 142 this article and describes possible future work direction.

II. RELATED WORK

143

144

As a dimensionality reduction technique, feature selection 145 algorithms try to find a lower-dimensional representation of 146 data via removing irrelevant features without altering the 147 original feature space [12]-[17]. Traditional feature selec- 148 tion methods can be grouped into three categories, that is, 149 filter, wrapper, and embedded approaches. Filter approaches 150

232

254

151 rank the features with correlation coefficients first and then 152 select the most suitable features. Peng et al. proposed a minimal redundancy and maximal relevance [18] algorithm, which 154 selects relevant features and simultaneously removes redun-155 dant features according to the mutual information. Another 156 filter method, fast correlation-based filter [19], exploits symmetrical uncertainty for feature selection. Wrapper approaches 158 apply a heuristic search strategy to determine feature subsets 159 and evaluate them based on the classification performance. 160 For example, Maldonado and Weber [20] proposed a wrap-161 per method for feature selection problems using support 162 vector machines (SVMs). The feature selection process and 163 classification model are separated and independent in the fil-164 ter and wrapper methods, and the wrapper methods might 165 suffer from high computational complexity especially for 166 high-dimensional data [21]. Embedded methods combine the 167 advantages of the filter and wrapper methods, which perform 168 the feature selection as part of its classicization process and obtain the feature subsets by optimizing the objective function. 170 such as a sparse Bayesian-based feature selection method [16]. 171 More recently, an evolutionary-computation-based [22] self-172 adaptive particle swarm optimization algorithm was proposed 173 for large-scale feature selection problem [23].

However, most of the traditional feature selection algo-175 rithms ignore the cause-effect relationships between features and the class attribute, and thus do not lend themselves to make predictions of the results of actions or interventions [24], [25]. To build better interpretability for data and robustness against 179 noise, causal feature selection methods were proposed, which 180 discover the MB of a target [1]. Pellet and Elisseeff [26] theoretically proved that MB is the optimal solution for feature 182 selection problem under the faithfulness condition. Therefore, 183 causal feature selection methods based on MB have attracted 184 more and more attention in recent years.

174

The first MB discovery algorithm for feature selection is the 186 KS [6] algorithm proposed by Margaritis and Thrun [5] KS discovers MBs by minimizing the cross-entropy loss without 188 theoretical guarantees to soundness. The GS [5] was the first 189 sound MB discovery algorithm, whose framework with the 190 growing phase and shrinking phase has become the basic strat-191 egy for the following algorithms. Tsamardinos et al. proposed 192 IAMB [7] to improve the GS by reordering the variables 193 in each iteration, which significantly improves the accuracy. 194 Based on IAMB, many variants have been developed, includ-195 ing inter-IAMB [7], fast-IAMB [27], and KIAMB [10]. These 196 algorithms are time efficient but require the number of samples to be exponential to the size of the MB, which means that insufficient samples will result in the performance degradation.

To improve the data efficiency while maintaining a reason-200 able time cost, a divide-and-conquer strategy for MB discovery proposed, that is, first finding the PC of a target, then 202 identifying spouses of the target. The MMMB [8] adopts the 203 divide-and-conquer strategy to search MBs, in which the data requirement is dependent on the topological structure rather 205 than the size of a variable set. Another early method, HITON-206 MB [9], interweaves the growing phase and the shrinking 207 phase in the PC discovery process, so that the false PC vari-208 ables can be excluded as early as possible. Pena et al. pointed out some errors in the PC discovery in the MMMB and 209 HITON-MB and, then, they added the double check strategy 210 to the MMMB framework and presented the PCMB [10] algo- 211 rithm. Based on PCMB, iterative parent-child-based search of 212 MB (IPCMB) [28] improves the time efficiency by connecting 213 target to all other variables and removing the false variables 214 in each iteration. De Morais and Aussem [11] proposed the 215 MBOR, which uses a weak MB learner (a fast but data- 216 inefficiency algorithm) to obtain the initial MB first and then 217 corrects the MB through a divide-and-conquer search, which 218 further improves the accuracy and data efficiency of the MB 219 discovery. Recently, Gao and Ji [29] discovered the coex- 220 istence of the spouses and the false parent-child variables, 221 and proposed a relatively efficient algorithm, simultaneous 222 MB (STMB), which improves the time efficiency of the MB 223 discovery.

Although existing causal feature selection algorithms 225 improve the data efficiency and accuracy, there are still sev- 226 eral true positives that cannot be identified, especially, with 227 the noise or small-sized samples [12], [30]. In this article, we 228 will focus on a type of incorrect CI tests occurred in the MB 229 discovery, and further improve the accuracy of MB discovery 230 through tackling this problem.

III. NOTATIONS AND DEFINITIONS

In this article, the capital letters (such as X, Y) represent 233 the random variables and the lowercase letters (such as x, y) 234 represent their values, the capital bold italic letters (such as 235 U, Z) denote variable sets. Specifically, let T denote the target 236 variable, and U denote the (discrete random) variable set.

Definition 1 (CI): Variables X and Y are conditionally independent given a variable set **Z** if $P(X, Y|\mathbf{Z}) = P(X|\mathbf{Z})P(Y|\mathbf{Z})$, 239 denoting as $X \perp Y|Z$. Similarly, $X \not\perp Y|Z$ represents that X and Y are conditionally dependent given a variable set Z.

Existing MB discovery algorithms use the G^2 -test [2] to 242 implement the CI test. In this article, we use the symbol 243 $dep(X, Y|\mathbf{Z})$ to represent the degree of the dependence between 244 X and Y conditioned on Z.

Definition 2 (Bayesian Network) [1]: Let \mathbb{P} denote the joint 246 probability distribution over a variable set U of a DAG \mathbb{G} . 247 The triplet $\langle U, \mathbb{G}, \mathbb{P} \rangle$ constitutes a BN, if $\langle U, \mathbb{G}, \mathbb{P} \rangle$ satisfies 248 the Markov condition: every variable is independent of any 249 subset including its nondescendant variables given its parents 250 in \mathbb{G} . In $\langle U, \mathbb{G}, \mathbb{P} \rangle$, the joint probability \mathbb{P} can be decomposed 251 into a product of conditional probabilities as follows:

$$P(U) = \prod_{X \in U} P(X|Pa(X))$$
 253

in which Pa(X) denotes the parents of X.

Some terms in the BN need to be declared here. If there 255 exists an edge from a variable (or node) X to Y, for example, 256 $X \to Y$, then X is a parent of Y and Y is a child of X. A 257 variable X is a spouse of Y if they share common child. In $_{258}$ this article, we denote by PC(X), the set of parent-child vari- 259 ables of X, and SP(X), the set of spouse variables of X. For 260 convenience, we will abbreviate parent-child and spouse as 261 PC and SP, sometimes. Based on the definition of BN, some 262 basic definitions in BN will be presented in the following.

Definition 3 (Blocked Path) [1]: A path π from variable A to B is blocked by a variable set \mathbf{Z} iff: 1) π contains a chain $A \to X \to B$ or $A \leftarrow X \to B$ with the middle variable $X \in \mathbf{Z}$ and 2) π contains a collider $A \to X \leftarrow B$ with $X \notin \mathbf{Z}$.

Definition 4 (d-Separation) [1]: In a DAG \mathbb{G} , variable set $\mathbf{Z} \subset U$ d-separates variables X and Y iff \mathbf{Z} blocks every path from X to Y, denoting as d-sep $(X, Y|\mathbf{Z})$.

With Definitions 3 and 4, we give the definition of faithful-272 ness condition.

Definition 5 (Faithfulness) [2]: Given a BN $\langle U, \mathbb{G}, \mathbb{P} \rangle$, \mathbb{P} is faithful to \mathbb{G} when for any $X, Y \in U$ and $\mathbf{Z} \subseteq U - \{X, Y\}$, 275 $X \perp Y | \mathbf{Z}$ in \mathbb{P} iff d-sep $(X, Y | \mathbf{Z})$ in \mathbb{G} .

Definition 5 shows that CI and d-separation are equivalent if the dataset and its underlying BN are faithful to each other. Thus, we have Theorem 1 as follows.

Theorem 1: In BN $\langle U, \mathbb{G}, \mathbb{P} \rangle$, for $X, T \in U$, there is an edge between X and T iff $X \not \perp T | \mathbb{Z}$ for $\forall \mathbb{Z} \subseteq U - \{X, T\}$.

Theorem 1 illustrates that if X is a PC variable of T, X and T are conditionally dependent for $\forall \mathbf{Z} \subseteq U - \{X, T\}$. Theorem 1 will help us to design the algorithm to discover the PC variables. With Definition 5, we give the definition (Definition 6) and property (Theorem 2) of MB in a faithful BN.

Definition 6 (Markov Boundary) [1]: In a faithful BN $\langle U, \mathbb{G}, \mathbb{P} \rangle$, the MB of a target variable T in \mathbb{G} is unique and consists of its parents, children, and spouses.

Theorem 2 [1]: Given the MB(T), X is independent of T 290 for any $X \in U - MB(T) - \{T\}$, that is, $X \perp T | MB(T)$.

According to Definition 6, for each variable, its MB can be easily "read" from the structure of the corresponding faithful BN. To understand the intuition in the perspective of causal learning, we consider that the MB includes the direct causes (parents), direct effects (children), and other direct causes of direct effects (spouses) of the class attribute [24].

Theorem 3 [3], [4]: The MB is the optimal solution for the feature selection problem.

Theorem 3 presents the significance of MB research, which confirms that we can transfer the feature selection problem into the MB discovery of the class attribute in a faithful BN.

IV. PCMASKING: TYPE OF INCORRECT CI-TESTS IN MARKOV BOUNDARY DISCOVERY

In this section, we focus on analyzing the three following questions: 1) Which variables are false negatives? 2) Why
ing questions: 1) Which variables are false negatives? 2) Why
can be are these variables discarded? and 3) Can we seek a theoretical solution to support the improvement of the algorithm?
To answer these questions, we first give two examples of
the PCMasking phenomenon using the existing algorithms in
Section IV-A, and then analyze the mechanism behind the
PCMasking phenomenon in Section IV-B. And finally, we anaIze lyze the effect of the phenomenon on existing causal feature
selection algorithms in Section IV-C.

314 A. Motivation

302

By running existing causal feature selection algorithms (PCMB, MBOR, and STMB) using the benchmark BN datasets with 1000 samples, we analyzed the incorrect CI tests

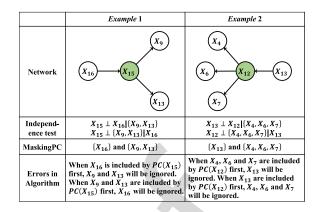


Fig. 2. Examples: two subnetworks of Alarm BN to demonstrate the incorrect CI tests using the existing algorithms. The figure shows the corresponding DAG with targets highlighted in green, the errors in the CI tests, the MaskingPC, and the errors in the algorithm caused by PCMasking.

in the discovered PC sets. We found that most of the false negatives (undetected PC variables) are independent of the target conditioning on other variables which have a strong correlation with the target (e.g., other PC variables). To further demonstrate this phenomenon, we take two subnetworks in the benchmark Alarm BN [31] as examples as follows.

Consider Example 1 in Fig. 2 and take X_{15} as the target, then X_{16} , X_{9} , and X_{13} are PC variables. According to 325 Theorem 1, $X_{16} \not = X_{15} | \{X_9, X_{13}\}$ and $\{X_9, X_{13}\} \not = X_{15} | X_{16}\}$ 326 hold. However, in our experiments, we have found that the 327 target X_{15} is independent of its parent X_{16} conditioning on its 328 children $\{X_9, X_{13}\}$, and the target X_{15} is independent of its 329 children $\{X_9, X_{13}\}$ conditioning on its parent X_{16} .

Using the other benchmark BN datasets, we also find many similar phenomenons in the experiments. Therefore, if we do 332 not tackle the type of incorrect CI tests, many true PC variables 333 may be discarded, leading to a low true positive discovery 324 accuracy. Furthermore, since the second type of causal feature 335 selection algorithms identifies the PC variables first, and then 336 finds the SP variables, this type of incorrect CI tests will lead 337 to the cascading errors during the SP discovery. Thus, it is 338 important to address the problem to improve the true positive 339 discovery accuracy. Before addressing this phenomenon, we 340 first give a formal definition as follows.

Definition 7 (PCMasking): In a variable set U, let PC(T) 342 denote the parent-child set for the target T. PC_{S1} and PC_{S2} 343 are the subsets of PC(T), and $PC_{S1} \cap PC_{S2} = \emptyset$. PC_{S1} and 344 PC_{S2} are PCMasking for T if the following conditions hold: 345

$$T \perp PC_{S1}|PC_{S2}, T \perp PC_{S2}|PC_{S1}.$$

346

We call PC_{S1} and PC_{S2} as Masking PC_{S3} .

Note that PC_{S1} and PC_{S2} in Definition 7 can not only be 348 single variables but also be sets with several variables. The 349 definition of PCMasking describes the results of the CI tests 350 rather than the mechanism shown in the BN, that is to say, if 351 the CI tests of an MB algorithm obtain the above independent 352 relationship, then there exists the PCMasking phenomenon. 363 For example, in Fig. 2, $\{X_{16}\}$ and $\{X_{9}, X_{13}\}$ are PCMasking 354 for X_{15} , and $\{X_{13}\}$ and $\{X_{4}, X_{6}, X_{7}\}$ are PCMasking for X_{12} . 355

356 B. Mechanism Analysis

This section focuses on explaining the mechanism of PCMasking. The direct reason for PCMasking is the complete dependence between variables which can be mathematically described as P(X=x|Y=y)=1. Obviously, complete dependence is a sufficient condition for PCMasking, which can be intuitively understood from an example. For Example 1 in Fig. 2, $P(X_{15}=0|X_{16}=0)=1.0$ and $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, then $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, and $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, then $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, then $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, then $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, and $P(X_{15}=0,X_{13}=1|X_{15}=0)=1.0$, then $P(X_{15}=0,X_{15}=1|X_{15}=0)=1.0$, then $P(X_{15}=0,X_{15}=1|X_{15}=0)=$

Theorem 4: In a dataset with variable set U, variable set S set S and S set S are S are the subsets of S set S and S set S are S such that S such

Proof: First, we prove the sufficiency of the condition. Since $\exists x, y, t \text{ s.t. } P(X = x | T = t) = P(Y = y | T = t) = 1$, then $\exists x \in X \perp T \text{ and } Y \perp T \text{ by Definition 1. Due to the complete dependence between } X \text{ and } T, T \text{ is independent with any otherwise given a variable set } X. \text{ Therefore, } Y \perp T | X \text{. And } X \perp T | Y \text{.}$ and be proved in the same way. Consequently, $X \text{ and } Y \text{ are PCMasking for } T. \text{. Similarly, we can prove that if } \exists x, y, t \text{ s.t.}$ and $Y \perp T | Y \text{.}$ and $Y \perp T | Y \text{.}$

Second, we utilize the concept of entropy in information theory to prove the necessity of the condition. Let H(X) denote the information entropy of X. Since $X \perp T|Y$ and $Y \perp T|X$, then the conditional information entropy can be simplified as follows:

$$\begin{cases}
H(X, T|Y) = H(X|Y) + H(T|Y) \\
H(Y, T|X) = H(Y|X) + H(T|X).
\end{cases}$$
(1)

388 According to the additivity of the information entropy, we have

$$\begin{cases}
H(X, Y, T) - H(Y) = H(X, Y) - H(Y) + H(T|Y) \\
H(X, Y, T) - H(X) = H(X, Y) - H(X) + H(T|X).
\end{cases} (2)$$

By solving the simultaneous equations in (2), we immediately obtain

$$H(T|X) = H(T|Y). \tag{3}$$

393 Obviously, there exist two cases that satisfy (3).

392

395

396

397

398

399

400

402

403

404

1) $H(T|X) = H(T|Y) \neq 0$: Since information entropy is based on a complete probability distribution, we suppose that $P(T = t_i | X = x_j) = p_{ij}$, $P(T = t_i | Y = y_k) = q_{ik}$, and $P(T = t_i | X = x_j, Y = y_k) = h_{ijk}$ ($i \in [1, |T|], j \in [1, |X|], k \in [1, |Y|]$), where |X| denotes the domain of X. We try to solve h_{ijk} via equation set

$$\begin{cases} \sum_{k=1}^{|Y|} h_{ijk} = p_{ij} \\ \sum_{j=1}^{|X|} h_{ijk} = q_{ik} \end{cases}$$
 (4)

where p_{**} and q_{**} are used as the parameters. Note that there are (|T-1|)(|X|+|Y|) equations and (|T-1|)(|X|)(|Y|) dependent variables. Therefore, the solution of (4) can be obtained iff H(T|X) = H(T|Y) = 0, contradicting the condition of the case. Consequently,

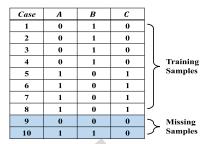


Fig. 3. Example of crucial samples: when the last two samples in the dataset are missing (highlighted in blue), the dependency between A and B (C) will change from an uncertain relationship to a deterministic relationship.

there exists no $P(T = t_i | X)$ and $P(T = t_i | Y)$ satisfying 406 the condition. As a result, case 1) does not hold true.

2) H(T|X) = H(T|Y) = 0: Since the condition of H(X) = 408 0 is that there exists x s.t. P(X = x) = 1, then we 409 can conclude that $\exists x, y, t$ s.t. P(T = t|X = x) = 1 and 410 P(T = t|Y = y) = 1. According to the symmetry of 411 the complete dependence, we can further conclude that 412 P(X = x|T = t) = 1 and P(Y = y|T = t) = 1.

Theorem 4 shows that if there exists PCMasking phenomenon, then it must be due to the complete dependence the between variables. Note that standard BN datasets (e.g., 417 Alarm) all satisfy the faithfulness condition, which makes the complete dependence not allowed in the original probability distribution in BN. The question will be, why does PCMasking phenomenon exist in the standard BN datasets?

Actually, the direct factor in determining the correctness of 422 CI test is the underlying probability distribution in the sam- 423 ples, instead of the original probability distribution in the BN, 424 we need to focus on the influence from the samples. Note 425 that the size of samples used in our experiments is 1000. 426 Although it is not small for the Alarm BN with 37 variables, 427 the lack of some crucial samples may cause a strong change 428 in the probability distribution of the variables, which makes 429 the datasets not satisfy the original distribution in the DAG. 430 The most significant change is that the dependency between 431 variables changes from an uncertain relationship to a deter- 432 ministic relationship. For example, suppose A, B, and C each 433 have space $\{0, 1\}$, and we have these samples in Fig. 3. We 434 can conclude from Fig. 3 that A and B (and C) are depen- 435 dent and the relationship between them is stochastic, that is, 436 given the value of A, we cannot determine the value of B or C. 437 However, when we suppose our sample contains the first eight 438 data items and the cases 9 and 10 (highlighted in blue) have 439 been lost, then variable dependency will become deterministic, 440 that is, P(B = 1|A = 0) = 1 and P(C = 1|B = 0) = 1.

Due to the loss of some crucial samples, the complete dependence appears in the underlying probability distribution of the samples, leading to the PCMasking phenomenon. 444 Obviously, it will further interfere with the MB discovery as 445 we discussed above. Unfortunately, there is no efficient way to 446 test the complete dependence for MB discovery. Therefore, we 447 need to exploit a variable-relationship-based method to detect 448 the MaskingPCs, instead of a dependence-test-based way. 449

450 C. Influence and Solution

According to the examples analyzed above, the PCMasking phenomenon will interfere with the recognition of PC variables, making existing methods discover false PC and SP sets. Given the problem of existing algorithms, we propose Theorem 5 to formalize the influence of PCMasking phenomenon on MB discovery and give an insight on how to detect the MaskingPCs in BN.

Theorem 5: In variable set U, let $PC_R(A) \subset U$ denote the real parent—child variable set of variable A, and $PC_O(A) \subset U$ denote the parent—child variable set of variable A found by the MB algorithm Ω . If a variable $Y \in PC_R(T)$ and a variable set $X \subseteq PC_R(Y)$ are PCMasking for T, then: 1) T might be ignored by Ω , that is, there exists the situation that $T \notin PC_O(Y)$ and $T \in PC_O(T)$.

Proof: We first prove that there exists the situation that $T \notin PC_O(Y)$. X and Y are PCMasking for T, we have $T \perp Y | X$ and Y based on Definition 7. Since $X \subseteq PC_R(Y)$, when the variables in X are all selected into the $PC_O(Y)$ before T, the existing algorithms will misjudge T as a non-PC variable of Y according to Theorem 1. Therefore, we prove that $T \notin PC_O(Y)$. Second, we prove $Y \in PC_O(T)$. Whatever order the variables in $PC_R(T)$ are selected to $PC_O(T)$, the variable Y will be added to $PC_O(T)$ and will not be deleted according to Theorem 1. Therefore, $Y \in PC_O(T)$.

In Theorem 5, Proposition 1) shows that existing algo-475 476 rithms do not guarantee the correct output in the BN with 477 PCMasking phenomenon. If the PC subset of target T is a 478 MaskingPC, then existing algorithms may fail to discover the 479 PC set of T under the influence of false CI tests. Taking 480 Fig. 2 as an example. Due to the PCMasking, we have ⁴⁸¹ $X_{16} \perp X_{15} | \{X_9, X_{13}\}$. According to Theorem 1 (criteria for 482 identifying PC variables), when variables X_9 and X_{13} are 483 selected into $PC_O(T)$ in advance, X_{16} will not be included 484 by $PC_O(T)$ since $T \perp X_{16} | \{X_9, X_{13}\}$ although $X_{16} \in PC_R(T)$. 485 Proposition 2) describes that the MaskingPC will break the 486 symmetry between the PC variables. Also taking Fig. 2 as an 487 example. Since there is no PCMasking phenomenon, then we 488 have $X_{15} \in PC_O(X_{16})$, while $X_{16} \notin PC_O(X_{15})$ according to the above analysis. Therefore, the symmetry between X_{15} and X_{16} 490 is broken in the output of the existing algorithms. Proposition 491 2) shows the possibility that we can detect the MaskingPCs 492 and simultaneously recover the discarded PC variables via the 493 symmetry property, which would be used for the improvement 494 of MB discovery algorithm.

V. CCMB AND PM-CCMB ALGORITHMS

This section presents the proposed MB discovery algorithms, CCMB in Section V-A and PM-CCMB in Section V-B.

498 A. CCMB

This section focuses on the specific design of the CCMB algorithm, including an innovative PC discovery process, in which there exists the handling of PCMasking phenomenon. The pseudocode of the CCMB is shown as Algorithm 1, which consists of three parts. In part 1 (line 3), CCMB discovers the PC variables using a subroutine called *FindPC* (detailed in

Algorithm 1 CCMB(T): Discover the MB of T

```
1: Input: Target variable T, variables set U.
2: Initialize the PC variable set, SP variable set and
    PCMasking table PC, SP, PCMTab \leftarrow \emptyset.
    {Part 1: Discover the PC variables.}
 3: PC \leftarrow FindPC(T)
    {Part 2: Detect the MaskingPCs and eliminate the effect
    of them.}
 4: for each X \in U - \{T\} do
      if T \in FindPC(X) and X \notin PC, then
5:
         PCMTab \leftarrow PCMTab \cup \{[T, X]\}
      end if
 7:
 8: end for
   for each [T, X] \in PCMTab do
      PC = PC \cup \{X\}.
11: end for
    {Part 3: Discover the spouse variables.}
12: for each Y \in PC do
13:
      for each X \in PC_Y do
         if X \notin PC, then
14:
             find Z s.t. T \perp X | \mathbf{Z} and T, X \notin \mathbf{Z}.
15:
             if T \perp X \mid \mathbf{Z} \cup \{Y\}, then
16:
                SP \leftarrow SP \cup \{X\}.
17:
            end if
18:
          end if
19:
      end for
20:
21: end for
22: Output: The Markov boundary of T, MB \leftarrow PC \cup SP.
```

Algorithm 2). *FindPC* will find all the true positives except 505 the PC variables with PCMasking phenomenon. Part 2 (lines 506 4–11) in CCMB detects the MaskingPCs and recovers the discarded variables. Based on the discovered PC sets in parts 1 508 and 2, CCMB calls part 3 (lines 12–21) to discover the SP set. 509

CCMB shares the same basic hypotheses with existing algorithms, that is, faithfulness and causal sufficiency. They will 511 be used for each theorem and its proof without restatement. 512 In the following text, we will describe the process of CCMB 513 step by step, and give some theoretical analyses. 514

The **part 1** in CCMB (Algorithm 1, line 3) identifies the PC variables except the MaskingPCs, whose process is detailed in 516 Algorithm 2. *FindPC* consists of three steps in an iteration. 517

Step 1 (Lines 4–14 of Algorithm 2): Build a candidate PC 518 set CanPC. Different from the existing methods, we exploit 519 a two-phase process to remove the false PC variables from 520 the current candidate PC set CanPC, which could effectively 521 reduce the number of iteration. The method is extracted from 522 Theorem 1: given variable set Z, if X and T are conditionally 523 independent, then X will be removed from the CanPC. Phase I 524 (lines 4–9) traverses the Z from current SPC. In line 5, Phase I 525 finds a conditioning set Sep[X] for each variable X, which can 526 minimize the conditional dependence between X and T, if X 527 and T are independent conditioned on Sep[X], then X will be 528 removed from the CanPC (lines 6–8). Phase II (lines 10–14) 529 selects the Z out of the SPC and specially considers a general 530 d-separation in BN, that is, the dependency is blocked by a 531

593

Algorithm 2 FindPC(T): Search the PC Subset of T

```
1: Input: Target variable T, variables set U.
 2: Initialize the PC variable subset SPC \leftarrow \varnothing, and the
    candidate PC variable set CanPC \leftarrow U - \{T\}.
 3: while CanPC \neq \emptyset do
       {Step 1: Find the candidate PC variables.}
       for each X \in CanPC do
 4:
 5:
          Sep[X] = arg min_{\mathbf{Z} \subseteq SPC} dep(T, X | \mathbf{Z}).
          if T \perp X | Sep[X], then
 6:
 7:
             CanPC \leftarrow CanPC - \{X\}.
          end if
 8:
       end for
 9:
       for each X, Y \in CanPC do
10:
          if X \perp Y and T \perp X|Y, then
11:
             CanPC \leftarrow CanPC - \{X\}.
12:
13:
          end if
       end for
14:
       {Step 2: Score the candidates and select the best.}
       for each X \in CanPC do
15:
          Score[X] = dep(T, X | Sep[X])
16:
       end for
17:
       Y = \arg\max_{X \in CanPC} Score[X].
18:
       SPC \leftarrow SPC \cup \{Y\}, CanPC \leftarrow CanPC - \{Y\}.
19:
       {Step 3: Delete the false variables.}
       for each X \in SPC do
20:
          if \exists Z \subseteq SPC - \{X\} s.t. T \perp X|Z, then
21:
             SPC \leftarrow SPC - \{X\}.
22:
          end if
23:
       end for
24:
25: end while
26: Output: The PC subset SPC.
```

532 single variable, which is also dependent on the two variables participating in CI test. In lines 11–13. Phase II removes variable X from the **CanPC** if X and Y are dependent on each other and Y blocks the dependency between T and X.

Theorem 6: Any true PC variable $X \in PC_R(T)$ is included 536 537 in CanPC until it is added to SPC.

Proof: According to Theorem 1, if X is a PC variable of 539 T, then $T \perp X|Z$ given $\forall Z \subseteq U - \{X, T\}$. Let us assume that variable $X \in PC_R(T)$, and X is removed from CanPC at lines ₅₄₁ 4–9 or lines 10–14 in Algorithm 2, then $\exists Sep[X]$ such that 542 $T \perp X | Sep[X]$. Thus, we have $X \notin PC_R(T)$, contradicting the $_{543}$ assumption. Therefore, any *X* ∈ $PC_R(T)$ will not be removed 544 from *CanPC* until it is added to *SPC*.

Step 2 (Lines 15-19 of Algorithm 2): Use scoring function 546 Score[X] to estimate the variables in CanPC. Score[X] is the $_{547}$ minimum of the conditional dependence between X and T, which has been calculated in line 5 of step 1. The dependence between X and T conditioned on Sep[X] can be used to score 550 the candidates in *CanPC* (line 16). We would select the variable with the highest score and add it to the SPC (lines 18 552 and 19).

Step 3 (Lines 20-25 of Algorithm 2): Detect the false vari-554 ables. As a heuristic method, the process above might intro-555 duce some false positives into the SPC. Therefore, lines 20–24

will remove the false variables from the current SPC accord- 556 ing to Theorem 1, that is, if there exists $Z \subseteq SPC - \{X\}$ such 557 that $T \perp X | \mathbf{Z}$, then X is a false positive.

Theorem 7: For variable $X \in U - \{T\}$: 1) if $X \in PC_R(T)$ 559 and there exists a variable set $S \subseteq PC_R(T)$ such that X and 560 **S** are PCMasking for T, then $X \in SPC$ or $X \notin SPC$; 2) if 561 $X \in PC_R(T)$ and there is no MaskingPCs for T, then $X \in SPC$; 562 and 3) if $X \notin PC_R(T)$, then $X \notin SPC$.

Proof: According to 1) of Theorem 5, proposition 1) of 564 Theorem 7 is true. In proposition 2), since there is no 565 MaskingPCs, if $X \in PC_R(T)$, then $T \perp X|Z$ with $\forall Z \subseteq 566$ $SPC - \{X\}$ according to Theorem 1. By Theorem 6, X will not 567 be removed until it is added to SPC. Therefore, proposition 2) 568 of Theorem 7 is true. In proposition 3), since $X \notin PC_R(T)$, then 569 $\exists \mathbf{Z} \subseteq SPC - \{X\}$ such that $T \perp X \mid \mathbf{Z}$. According to Theorem 6, 570 proposition 3) is also true.

The Part 2 in CCMB (Algorithm 1, lines 4–11) is the 572 biggest improvement compared with other methods. Based 573 on the output of FindPC (Algorithm 2), CCMB performs the 574 cross-check and complement processes, which can shield the 575 PCMasking phenomenon and simultaneously find more true 576 positives. The cross-check process (lines 4–8) aims to detect 577 the MaskingPCs and record the variables with its correspond- 578 ing target into the PCMasking table PCMTab. This process 579 utilizes the property that the MaskingPCs could break the sym- 580 metry between the PC variables and the target (proposed in 581 Theorem 5). Specifically, if T is included in the PC set of 582 X while the PC set of T does not include X, then the cross- 583check process can detect the MaskingPCs by the asymmetry 584 between X and T. After that, the complement process (lines 585) 9-11) repairs the asymmetry and completes the PC set.

The **Part 3** in CCMB (Algorithm 1, lines 12–21) discovers 587 the SP set. It uses the topology information of the BN to find 588 the colliders in the PC set of T. The spouses are selected from $_{589}$ the union of the PC sets of the PC variables of T (PC_Y in line 590 13 denotes the PC set of variable Y found by CCMB). To 591 illustrate the correctness of CCMB, Theorem 8 is proposed 592 and proved as follows.

Theorem 8: CCMB outputs the correct MB.

Proof: First, we prove that CCMB can output the correct 595 PC set. Based on proposition 3) of Theorem 7, CCMB can 596 delete all false PC variables. According to proposition 1) of 597 Theorem 7, some of the PCMasking variables may not be 598 included in **SPC**. Denoting one of them as X, then $T \in PC_O(X)$ 599 according to proposition 2) of Theorem 4. Therefore, the 600 PCMasking variables can also be correctly found by lines 4–11 601 of Algorithm 1. The remaining PC variables can be found 602 according to proposition 2) of Theorem 7. Second, we prove 603 that CCMB can find the correct SP set. For a variable $X \in SP$, 604 X has a common child with T, which can be found by lines 605 12–21 in Algorithm 1. Therefore, CCMB can find the correct 606 SP set. Summarizing, CCMB outputs the correct MB.

To further explain the algorithm, we will take Fig. 2 as 608 an example and revisit the steps of CCMB that were given 609 in Algorithms 1 and 2. Initially, $CanPC = \{X_{16}, X_9, X_{13}\}$ in 610 the first iteration and there is no variable removed from the 611 CanPC in lines 6-8 of Algorithm 2. Let us suppose a spe- 612 cial case to show the effect of detecting the MaskingPCs, 613

IEEE TRANSACTIONS ON CYBERNETICS

614 that is, $dep(X_9, X_{15}) > dep(X_{13}, X_{15}) > dep(X_{16}, X_{15})$ and 615 $\operatorname{dep}(X_{13}, X_{15} | \{X_9\}) > \operatorname{dep}(X_{16}, X_{15} | \{X_9\})$. Then, variable X_{16} 616 will be selected into the **SPC** in the first iteration, and X_9 and 617 X_{13} will be removed from the *CanPC* according to the lines 11-13 due to the PCMasking phenomenon as mentioned in 619 Section IV. We conclude that the outputs of Algorithm 2 are $SPC[X_{15}] = \{X_{16}\} \text{ and } SPC[X_9] = \{X_{15}\}, SPC[X_{13}] = \{X_{15}\}.$ Note that, existing algorithms will take three iterations to finish 622 the PC discovery and obtain an incorrect result, while CCMB 623 discovers the PC efficiently in only one iteration and con-624 tinues to select the undetected MaskingPCs. In Algorithm 1, 625 lines 4–8 traverse the outputs of Algorithm 2 and build the

8

631 B. Pipeline Machine: Acceleration Strategy

As the second type of causal feature selection algorithms, 633 CCMB has better accuracy but lower efficiency. In this section, we will propose an acceleration strategy to guarantee that 635 CCMB can be efficient in relatively large-scale datasets.

626 **PCMTab** = { $[X_{15}, X_9], [X_{15}, X_{13}]$ }, which means that the

symmetry between X_{15} and X_9 (X_{13}) are broken due to the

628 PCMasking phenomenon. Lines 9-11 repair the asymmetry and the outputs of Algorithm 1 are $SPC[T] = \{X_9, X_{13}, X_{16}\}.$ 630 Hence, our proposed CCMB obtains the expected outputs.

First, we point out the problems in the algorithms, and 636 637 take the pseudocodes in CCMB as examples to illustrate our idea. The MB discovery algorithms are designed based on the 639 CI tests, implemented by G^2 -test, which is a time-consuming 640 process. We found that there are a large number of iterative 641 processes consisting of the same CI tests, which are executed 642 multiple times by the algorithm. For example, when search- $_{643}$ ing for the PC set of variable X, we need to perform CI tests 644 between Y and X conditioned on different Z, while the same 645 tests will be performed when searching for PC set of variable Y. Lines 3–10 in Algorithm 1 are also the steps that need to be 647 performed only once between a pair of variables. The redun-648 dant tests exist not only between different variables, but also between different iterations of the same variable. Take lines 650 20–24 of Algorithm 2 as an example. SPC is similar between 651 adjacent iterations since only a small number of variables are 652 added or removed. Therefore, there are a large number of duplicate CI tests. In addition, redundant tests also exist in the 653 same iteration. For example lines 6 and 21 in Algorithm 2.

Inspired by the problems above, we propose an acceleration 656 methodology for the MB discovery algorithm called pipeline machine (PM). The main idea of PM is using a buffer to 658 organize the redundant CI test results involved in the MB dis-659 covery process, so that more efficient search can be exploited 660 to replace the complex calculations.

The structure of PM is shown in Fig. 4. PM is essentially a 661 662 linked list consisting of several variable cells (in the red dashed 663 frame). Each variable cell points to a CI test Information Table (in the blue dashed frame), which stores the new CI test results involved in the iteration of the corresponding variable being added to the MB. For the convenience of searching, the CI 667 test Information Table is organized into a two-level structure 668 in which CI tests with the same size of conditioning sets are

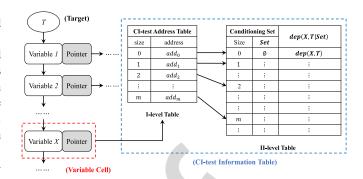


Fig. 4. PM for MB discovery algorithm. The skeletal structure is a linked list consisting of several variable cells. The red dashed part denotes the variable cell and the blue dashed part denotes the CI test Information Table which is used to store the results of the CI tests that have been executed.

stored in an II-level Table, and the addresses of the II-level 669 Table are stored in the I-level Table.

When a new variable is added to the MB, the PM adds a 671 variable cell to the end of the linked list. Before executing 672 a CI test, the PM-algorithm first determines whether the CI 673 test is already in the PM via the linked list structure of the 674 PM. Specifically, if all the variables in the conditioning set 675 (of the CI test) are in the linked list of the PM, then the CI 676 test is in the PM, and its location is in the CI-test Information 677 Table corresponding to the conditioning set variable closest 678 to the end of the linked list. For example, in the line 6 of 679 Algorithm 2, if all the variables in Sep[X] are in the linked 680 list of the PM, then the result of dep(X, T|Sep[X]) can be 681 obtained from the CI-test Information Table corresponding to 682 the variable in the Sep[X] closest to the end of the linked 683 list. As can be seen from the mechanism, PM reduces the 684 computational cost of redundant CI test, thus improving the 685 efficiency of the algorithm.

In this article, we refer CCMB with PM as PM-CCMB, indicating that it is an enhanced version based on PM. Note that 688 PM-CCMB could maintain the same accuracy with CCMB and 689 simultaneously improve the time efficiency through employing 690 the PM to store the repetitive computations, thus, PM-CCMB 691 need to require additional memory for the PM. We will discuss 692 the space complexity of PM-CCMB in the following. Assume 693 that the largest size of conditioning sets during the PC search 694 is m, then the size in the II-level Table (in Fig. 4) of variable 695 T is between 0 and max $\{m, |SPC_T|\}$, where SPC_T denotes the 696 current PC set after a certain iteration. Since there are $|PC_T|$ 697 II-level Tables, then there are $\sum_{i=0}^{|PC_T|-1} \sum_{j=0}^{\max\{m,i\}} \binom{i}{j}$ CI-test 698 results need to be stored. Therefore, PM-CCMB needs to store 699 $\sum_{T \in U} \sum_{i=0}^{|PC_T|-1} \sum_{j=0}^{\max\{m,i\}} {i \choose j}$ CI-test results.

VI. EXPERIMENTAL STUDIES

700

In this section, we present the experimental studies 702 of the proposed algorithms, CCMB and PM-CCMB. In 703 Section VI-A, we first use the subnetworks of Alarm BN in 704 Fig. 2 to demonstrate the effectiveness of CCMB to detect the 705 MaskingPCs. Then, we perform experiments on 12 standard 706 BN datasets in Section VI-B to evaluate several performance 707

708 aspects of the algorithms, including accuracy and time effi-709 ciency. In order to validate the performance in feature selection 710 problem, extensive experiments and statistical analysis are per-711 formed on a newly developed electroencephalography (EEG) dataset, SEED [32], [33], to compare our algorithms with other state-of-the-art methods in Section VI-C.

To illustrate the effectiveness of the proposed methods, the our state-of-the-art MB discovery algorithms are compared. 715

714

716

717

718

719

720

721

722

723

724

728

729

730

731

732

749

- 1) IAMB [7]: IAMB is the first type of causal feature selection algorithms, focusing on time efficiency.
- 2) PCMB [10]: PCMB is the second type of causal feature selection algorithms, which aims to improve accuracy.
- 3) MBOR [11]: MBOR combines the idea of two types of algorithms to improve accuracy.
- STMB [29]: STMB is a recently proposed algorithm, which incorporates the double check process of PCMB into the SP discovery to improve the time efficiency.

In addition, we also select two well-established information-726 theoretical-based feature selection algorithms as references in the experiment of emotion recognition task.

- 1) FCBF [19]: A fast correlation-based filter, which exploits the symmetrical uncertainty for feature selec-
- 2) mRMR [18]: An algorithm of removing redundant features while ensuring maximum correlation.

To measure the strength of the conditional dependence 733 734 between variables, all the MB discovery algorithms use the 735 G^2 -test to implement the CI tests as previous work ([11], [29], 736 etc.) at a significance level of 0.01. All codes are implemented 737 in C++, and all experiments are conducted on a computer 738 with Inter i5-8500 3.00-GHz CPU and 16-GB memory.

739 A. Alarm Subnetwork: The Effectiveness for Detecting 740 MaskingPCs

The proposed CCMB has been proved to detect the 742 MaskingPCs and shield its impact. This experiment selects the 743 two subnetworks of Alarm BN in Fig. 2 to demonstrate the effectiveness of CCMB to solve the PCMasking phenomenon. The two examples in Fig. 2 take variables X_{15} and X_{12} as the 746 target, respectively. To validate the existence of PCMasking 747 phenomenon, we test the CI with 1000 samples and obtain 748 the results as follows:

$$X_{15} \perp X_{16} | \{X_9, X_{13}\}, X_{15} \perp \{X_9, X_{13}\} | X_{16}$$

 $X_{12} \perp X_{13} | \{X_4, X_6, X_7\}, X_{12} \perp \{X_4, X_6, X_7\} | X_{13}$

751 which means that X_{16} and $\{X_9, X_{13}\}$ are PCMasking for X_{15} , and X_{13} and $\{X_4, X_6, X_7\}$ are PCMasking for X_{12} . According to 753 Theorem 4, the existing algorithms cannot find out the correct 754 PC sets of X_{15} and X_{12} . Table I provides the PC variables of 755 X_{15} and X_{12} found by our algorithm and other MB discovery 756 algorithms.

For Example 1 in Fig. 2, variables X_9 and X_{13} have 758 stronger correlation with variable X_{15} , such that they will 759 enter into $PC(X_{15})$ in advance and then prevent variable 760 X_{16} . When using $\{X_9, X_{13}\}$ as the conditioning set, we have $|\text{dep}(X_{15}, X_{16}|\{X_9, X_{13}\})| < -1.0$. Therefore, X_{16} is misjudged 762 as a non-PC variable according to Theorem 1. This error will

TABLE I SEARCHED PC VARIABLES ON ALARM SUBNETWORK WITH PCMASKING PHENOMENON

Algorithms	Results of Example 1	Results of Example 2
IAMB	$X_{16} \notin MB(X_{15})$	$X_{13} \notin MB(X_{12})$
	$X_{15} \in MB(X_{16})$	$X_{12} \in \mathbf{MB}(X_{13})$
PCMB	$PC(X_{15}) = \{X_9, X_{13}\}$	$PC(X_{12}) = \{X_4, X_6, X_7\}$
	$PC(X_{16}) = \emptyset$	$PC(X_{13}) = \emptyset$
MBOR	$PC(X_{15}) = \{X_9, X_{13}\}$	$PC(X_{12}) = \{X_4, X_6, X_7\}$
	$PC(X_{16}) = \{X_{15}\}$	$PC(X_{13}) = \{X_{12}\}$
STMB	$PC(X_{15}) = \{X_9, X_{13}\}$	$PC(X_{12}) = \{X_4, X_6, X_7\}$
SIMD	$PC(X_{16}) = \emptyset$	$PC(X_{13}) = \emptyset$
ССМВ	$PC(X_{15}) = \{X_9, X_{13}, X_{16}\}$	$PC(X_{12}) = \{X_4, X_6, X_7, X_{13}\}$
	$PC(X_{16}) = \{X_{15}\}$	$PC(X_{13}) = \{X_{12}\}$

STATISTICAL INFORMATION OF THE STANDARD BN DATASETS

Data set	Alarm	Alarm3	Child	Child3	Insurance	Insurance3
#Variables	37	111	20	60	27	81
#Edges	46	149	25	79	52	163
Data set	Alarm5	Alarm10	Child5	Child10	Insurance5	Insurance10
#Variables	185	370	100	200	135	270
#Edges	265	570	126	257	281	556

occur in all other existing algorithms. Moreover, STMB and 763 PCMB further misjudge that X_{15} is not a PC variable of X_{16} . In 764 contrast, our proposed CCMB achieves the complete PC sets 765 for variables X_{15} and X_{16} due to the cross-check and complement processes that can shield the influence of the PCMasking 767 for X_{15} . Similarly, when using $\{X_4, X_6, X_7\}$ as the conditioning 768 set in Example 2, X_{13} will be misjudged as a non-PC variable 769 in all the existing algorithms while CCMB can avoid errors 770 on MaskingPCs.

B. Standard BN Datasets: The Accuracy and Time Efficiency 772 for MB Discovery

In this section, we conduct experiments on standard BN 774 datasets [34] to evaluate the performance of CCMB and 775 other algorithms. The standard BN data contains 12 network 776 datasets. Table II provides the statistical information of these 7777 datasets, including the number of variables, the number of 778 edges, and the training size in the experiments. The experiments consist of two parts, to verify the accuracy and time 780 efficiency of the proposed algorithms, respectively.

We run all the MB algorithms for each variable and repeat 782 these algorithms 20 times with different samples. The size 783 of training samples turns from {500, 1000, ..., 4000}, respec- 784 tively. We compare the accuracy and time efficiency of the 785 algorithms under the same sets of samples and analyze the 786 changes of the performance with increased training sample 787

Accuracy: The frequently used metrics Distance [7], [10], 789 [11], [29] is adapted to measure the accuracy of MB vari- 790 ables searching. The Distance measures the distance between 791 the detected MB and the true MB, calculated by: Distance = 792 $\sqrt{(1 - \text{Precision})^2 + (1 - \text{Recall})^2}$, where the *Precision* is the 793 fraction of retrieved true positives over the total amount of 794 detected MB variables, and the Recall is the fraction of 795 retrieved true positives over the total amount of true MB vari- 796 ables. Thus, the lower Distance indicates the detected MB is 797 closer to the true MB.

IEEE TRANSACTIONS ON CYBERNETICS

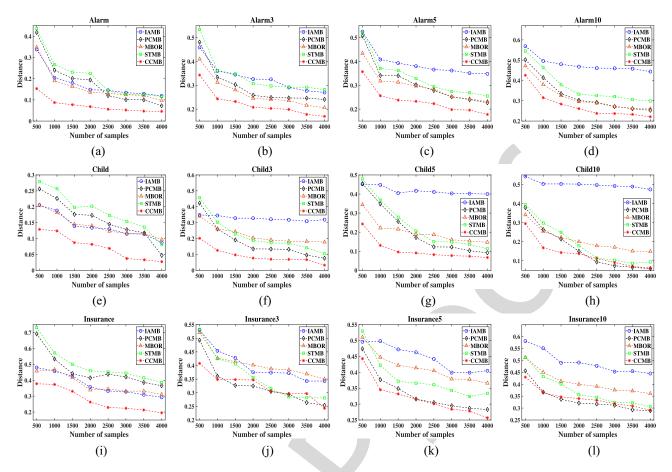


Fig. 5. Results of the MB discovery experiments for the accuracy of CCMB and other algorithms on the 12 standard BN datasets. Note that the curves of PM-CCMB and CCMB are identical since PM-CCMB only optimizes the algorithm structure of CCMB. Therefore, we do not make a distinction. (a) Alarm. (b) Alarm3. (c) Alarm5. (d) Alarm10. (e) Child3. (g) Child3. (g) Child40. (i) Insurance. (j) Insurance3. (k) Insurance5. (l) Insurance10.

Fig. 5 shows the average Distance variation curves of CCMB and other algorithms with respect to the number of samples. From Fig. 5, we can observe that all the algorithms trend to achieve a lower Distance with more samples, and CCMB consistently performs better than others with different scales of samples on 9 out of 12 datasets. Since CCMB considers the influence of PCMasking phenomenon on MB 806 discovery, more true positives are detected, which improves the recall, thereby, improves the accuracy of the algorithm. The 808 Distance of CCMB is similar with PCMB but also smaller than other algorithms on the Child 10, Insurance 3, and 810 Insurance 10 with large-scale samples (#Variables > 1000). This is because CCMB may introduce some false positives into 812 the MB while finding more true positives, which causes the 813 precision of CCMB to drop slightly, resulting in the accuracy 814 of CCMB similar to PCMB. We also note that CCMB consis-815 tently outperforms the other four MB discovery algorithms on datasets under small-scale samples (#Variables \le 1000), which demonstrates the significant superiority of CCMB in data efficiency. Especially, on the Child dataset, CCMB uses fewer samples (500 samples) to achieve a small Distance while other algorithms need 1500-2500 samples to achieve a similar Distance. With fewer samples, CCMB can find more accurate MBs, while other algorithms cause more errors due to 823 insufficient samples. When the sample size reaches a certain

scale, the existing algorithms can avoid some true positives 824 being ignored, making their accuracy close to our proposed 825 CCMB. Therefore, the significant superiority of CCMB under 826 small-scale samples reflects that our proposed CCMB is more 827 data-efficient. In general, CCMB can significantly improve the 828 accuracy in comparison to the state-of-the-art algorithms. 829

We recall that PM-CCMB only optimizes the executing process of CCMB. Therefore, in the above experiments, the curves of PM-CCMB and CCMB are identical so that we do not make a distinction. In the following text, we will see the superiority of PM-CCMB over CCMB and other state-of-the-art algorithms.

Time Efficiency: We recorded the CPU time for each dataset in the above experiments. Fig. 6 shows the logarithmic time 837 variation curves of CCMB, PM-CCMB, and other algorithms 838 with respect to the number of samples. From Fig. 6, we can 839 observe that the CPU time of CCMB is slightly higher than 840 PCMB and other algorithms, it is mainly because CCMB introduces the cross-check and complement processes to find more 842 true variables, which is a time-consuming step. To improve the 843 time efficiency, we proposed the PM. Therefore, PM-CCMB 844 consistently performs better than others with different scales 845 of samples on 10 out of 12 datasets, which demonstrates that 846 PM can significantly improve the computational efficiency 847 of CCMB through organizing and storing the intermediate 848

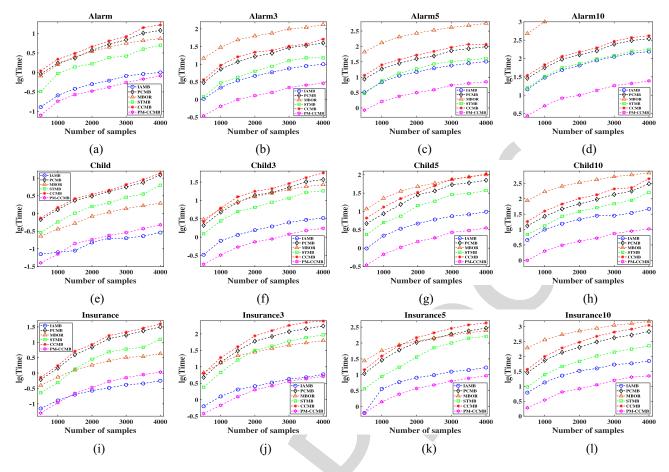


Fig. 6. Results of the MB discovery experiments for the time efficiency of CCMB, PM-CCMB, and other algorithms on the 12 standard BN datasets. (a) Alarm. (b) Alarm3. (c) Alarm5. (d) Alarm10. (e) Child (f). Child3. (g) Child5. (h) Child10. (i) Insurance. (j) Insurance3. (k) Insurance5. (l) Insurance10.

results. Especially, PM-CCMB is even faster than the most time-efficient IAMB. We also note that PM-CCMB is similar with or slightly higher than IAMB on the Child and Insurance. Mainly because, it will take more time to build the PM; then, the PM saves in some datasets with fewer variables (such as Child and Insurance). Even so, PM is still an effective accelerator for MB discovery algorithms. In general, PM-CCMB can significantly improve the time efficiency in spite of its impressive accuracy.

858 C. Emotion Recognition: The Effectiveness for Solving 859 Feature Selection Problem

This section employs the proposed MB discovery algorithms in feature selection task to solve the emotion recognition problem. A newly developed EEG dataset, SEED [32], will be used to evaluate the performance of PM-CCMB. The SEED dataset contains the EEG signals of 15 subjects, which are regarded as 15 datasets. Each subject watched 15 emotional film clips while the EEG signals were recorded by 62-channel symmetrical electrodes (shown in Fig. 7). The features are extracted from five common frequency bands, namely, Delta (1–3 Hz), Theta (4–7 Hz), Alpha (8–13 Hz), Beta (14–30 Hz), and Gamma (31–50 Hz), and each frequency band has 62-671 channel neural signatures. Therefore, there are 310 features in

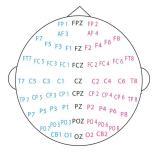


Fig. 7. Layout of 62 channel symmetrical electrodes on the EEG.

the SEED. The emotional labels (negative, neutral, and pos- 872 itive) of the film clips are used as the target. Each subject 873 performed the emotion experiments in three separate sessions 874 with an interval of about one week or longer, and each session 875 contains 3394 samples (1120 negative samples, 1104 neutral 876 samples, and 1170 positive samples).

In our experiments, the differential entropy (DE) features 878 are chosen for emotion recognition due to its better discrimination [32]. For each subject, we randomly choose 1000 samples 880 as the training set, and 500 samples as the test set. We compare the emotion recognition and feature selection effectiveness of 882 PM-CCMB with other algorithms on the test set.

We adopt three classifiers, that is, linear SVM, AdaBoost, 884 and Naive Bayes to compute their classification accuracies 885

Classifier	Subject	IAMB	PCMB	MBOR	STMB	FCBF	mRMR	РМ-ССМВ
	#1	87.06±5.30	92.87±3.98	87.87±6.22	90.58±4.31	97.37±2.62	92.07±1.25	99.91±0.22
	#2	67.04 ± 8.00	94.76 ± 4.47	67.39 ± 8.41	66.33 ± 7.25	96.98 ± 2.17	97.15 ± 0.93	99.99 ± 0.01
	#3	67.15 ± 8.44	96.09 ± 1.98	67.15 ± 8.44	96.32 ± 1.94	93.28 ± 6.78	97.89 ± 0.23	99.99 ± 0.01
	#4	69.01 ± 5.54	91.45 ± 3.42	69.18 ± 6.45	90.20 ± 1.97	98.56 ± 0.37	96.99 ± 1.36	100.00 ± 0.00
	#5	69.14 ± 2.73	96.14 ± 1.18	69.25 ± 2.49	95.54 ± 1.73	98.52 ± 1.10	96.65 ± 1.23	100.00 ± 0.00
	#6	72.36 ± 2.50	94.88 ± 3.75	72.36 ± 2.50	91.35 ± 4.62	99.50 ± 0.32	99.89 ± 0.11	100.00 ± 0.00
LibSVM	#7	77.09 ± 7.02	96.05 ± 1.26	80.25 ± 4.29	95.32 ± 2.38	98.98 ± 1.00	95.76 ± 0.43	99.99 ± 0.01
	#8	68.01 ± 17.14	94.04 ± 4.49	68.34 ± 17.68	90.52 ± 6.73	99.35 ± 0.61	98.35 ± 0.46	100.00 ± 0.00
	#9	70.55 ± 10.20	88.50 ± 4.86	65.79 ± 13.70	84.32 ± 3.74	82.28 ± 7.83	94.74 ± 1.15	99.90 ± 0.10
	#10	73.21 ± 4.71	87.86 ± 6.20	73.62 ± 4.70	85.69 ± 5.43	97.56 ± 3.19	97.96 ± 0.53	99.97 ± 0.02
	#11	82.74 ± 6.48	91.03 ± 8.41	83.14 ± 6.55	92.06 ± 7.42	99.77 ± 0.35	96.11 ± 1.13	$100.00\!\pm\!0.00$
	#12	65.24 ± 9.98	87.25 ± 3.41	66.35 ± 8.13	85.26 ± 4.74	98.74 ± 0.12	99.72 ± 0.23	$100.00\!\pm\!0.00$
	#13	54.85 ± 4.23	88.39 ± 6.08	55.27 ± 4.48	85.67 ± 5.32	99.28 ± 0.49	98.42 ± 0.59	100.00 ± 0.00
	#14	78.94 ± 9.25	81.02 ± 5.10	79.62 ± 10.77	83.44 ± 6.52	92.18 ± 0.97	92.68 ± 0.41	99.98 ± 0.01
	#15	58.48 ± 10.44	97.45 ± 3.19	58.70 ± 10.24	98.19 ± 1.14	99.70 ± 0.30	100.00 ± 0.00	100.00 ± 0.00
	average ranks	6.73	4.00	6.13	4.73	2.80	2.53	1.07
	#1	82.15±8.97	84.60±7.64	83.82±6.71	82.12±6.99	86.94 ± 7.25	90.10±1.18	98.08±1.26
	#2	58.08 ± 2.99	71.21 ± 5.82	57.94 ± 3.18	69.20 ± 5.97	82.22 ± 5.64	82.54 ± 0.76	90.77 ± 3.22
	#3	58.13 ± 3.84	72.04 ± 5.50	58.13 ± 3.84	73.06 ± 5.60	78.67 ± 1.37	88.37 ± 0.96	87.31 ± 0.95
	#4	65.53 ± 4.32	69.89 ± 4.34	66.33 ± 3.75	66.35 ± 4.50	90.80 ± 6.10	88.83 ± 1.35	91.98 ± 3.46
	#5	67.29 ± 2.70	77.49 ± 2.36	67.03 ± 2.97	78.25 ± 2.26	82.94 ± 8.71	81.77 ± 0.96	93.13 ± 1.00
	#6	75.41 ± 3.89	83.38 ± 7.38	75.41 ± 3.89	81.25±5.69	92.17±2.65	94.83 ± 2.62	95.40 ± 0.66
	#7	77.44 ± 7.22	90.08 ± 4.26	77.47 ± 7.19	88.25 ± 6.31	91.31 ± 0.37	96.15 ± 0.79	97.23 ± 0.53
Adaboost	#8	66.63 ± 15.76	77.30 ± 8.19	66.63 ± 15.76	76.46 ± 8.37	91.26 ± 2.65	93.69 ± 1.12	94.07 ± 2.04
	#9	67.96 ± 8.10	74.08 ± 7.20	66.34 ± 10.17	71.52 ± 8.12	79.12 ± 4.77	90.18 ± 2.25	94.66 ± 0.72
	#10	65.60 ± 2.22	67.83 ± 5.81	65.83 ± 2.56	63.26 ± 4.32	82.10 ± 2.96	91.14 ± 1.16	86.79 ± 0.97
	#11	78.99 ± 6.00	77.81 ± 4.88	79.11 ± 5.95	76.42 ± 5.23	91.56±1.28	90.53 ± 0.51	93.66 ± 0.99
	#12	64.15 ± 8.97	81.22 ± 6.35	67.32 ± 9.63	80.79 ± 4.53	88.74 ± 1.12	89.48 ± 0.92	91.26 ± 1.35
	#13	53.67 ± 4.04	75.53 ± 11.70	53.64 ± 3.26	72.56 ± 10.25	90.69 ± 4.51	94.56 ± 1.01	95.93 ± 1.23
	#14	61.72 ± 9.43	73.59 ± 5.32	61.28 ± 10.52	75.26 ± 6.48	91.25 ± 0.74	88.95 ± 2.74	96.73 ± 1.95
	#15	59.23 ± 13.37	95.45 ± 2.64	59.42 ± 13.19	96.42 ± 2.76	98.53 ± 0.29	97.51 ± 0.63	99.53 ± 0.46
	average ranks	6.20	4.40	6.27	5.13	2.67	2.20	1.13
Naive Bayes	#1	80.97±11.52	84.61±8.63	84.09±8.17	85.25±6.69	85.59±2.84	85.78±0.47	92.73 ± 2.83
	#2	54.54 ± 6.71	76.57 ± 8.54	54.67 ± 6.55	75.26 ± 9.94	80.37 ± 6.50	80.47 ± 0.62	85.96 ± 2.79
	#3	60.09 ± 5.22	76.82 ± 5.31	60.09 ± 5.23	74.36 ± 7.10	75.78 ± 1.61	85.37 ± 1.12	78.22 ± 3.41
	#4	69.95 ± 2.36	74.40 ± 3.22	70.61 ± 2.40	72.22 ± 4.85	92.74 ± 5.58	90.99 ± 1.22	92.68 ± 1.92
	#5	73.09 ± 3.90	88.81 ± 0.20	74.19 ± 3.43	85.12 ± 2.23	87.90 ± 5.28	90.13 ± 0.81	95.30 ± 0.96
	#6	73.76 ± 7.04	85.25 ± 9.90	73.76 ± 7.04	83.25 ± 8.02	95.83 ± 2.34	93.66 ± 1.95	96.77 ± 1.88
	#7	82.57 ± 4.76	89.60 ± 3.15	81.30 ± 5.87	84.59 ± 4.77	94.14 ± 1.19	91.24 ± 0.96	97.30 ± 0.23
	#8	72.97 ± 17.34	83.77 ± 6.15	73.39 ± 17.96	81.69 ± 3.26	94.06 ± 1.93	86.22 ± 0.42	100.00 ± 0.00
	#9	70.72 ± 7.36	79.06 ± 5.26	68.29 ± 9.56	80.15 ± 4.83	78.12 ± 1.86	81.99 ± 1.74	85.76 ± 1.62
	#10	69.57 ± 5.54	74.44 ± 7.30	69.66 ± 5.67	75.62 ± 5.65	85.77 ± 5.57	83.30 ± 0.56	91.01 ± 1.19
	#11	81.27 ± 5.08	78.22 ± 5.46	81.42 ± 5.10	77.42 ± 6.72	89.59 ± 2.02	81.74 ± 1.16	91.75 ± 1.08
	#12	65.39 ± 8.74	77.13 ± 4.32	66.25 ± 9.98	78.32 ± 5.65	87.36 ± 1.03	83.98 ± 1.99	91.28 ± 1.14
	#13	60.40 ± 4.48	83.09 ± 9.27	60.22 ± 4.27	84.15 ± 6.74	93.70 ± 4.88	92.56 ± 1.32	99.20 ± 0.56
	#14	71.33 ± 7.69	84.13 ± 4.56	73.25 ± 6.98	82.61 ± 3.42	87.35 ± 0.61	85.46 ± 2.85	89.15 ± 1.36
	#15	63.30 ± 10.71	97.16 ± 1.95	62.80 ± 11.33	96.99 ± 1.74	97.69 ± 1.23	98.68 ± 0.64	99.67 ± 0.71
	average ranks	6.47	4.27	6.27	4.73	2.60	2.53	1.13

achieved by using the selected feature subsets. In this experiment, the regularization parameter of linear SVM is tuned from $\{10^{-2}, \ldots, 10^2\}$ by the grid search strategy, and the number of ensemble learning cycles is 50 in the AdaBoost. We first run PM-CCMB and other algorithms to search its MB (or feature subset) on the training samples, then use the selected feature subset to train these classifiers on the training samples and, finally, test the accuracy on the test samples. The experiments are repeated 20 times with different training and test samples, and the average classification accuracy of different classifiers using the feature subsets selected by different algorithms are given in Table III.

12

From Table III, we note that the classifiers using the features selected by PM-CCMB achieve the highest or highly competitive classification accuracy on most of the subjects, which shows the effectiveness of PM-CCMB to select the relevant features. Specifically, the features selected by PM-CCMB can train a completely correct SVM classifier, while any other algorithms cannot achieve the accuracy. Due to the influence of the PCMasking phenomenon, the dependence between the target and its relevant features may be blocked by other relevant features. Therefore, existing MB discovery algorithms may ignore some of the relevant features resulting in the low accuracy of classification. By detecting the MaskingPCs, our algorithm can find more relevant features and significantly improve the accuracy of features selection. Through comparing with two well-established feature selection algorithms, we can conclude that PM-CCMB can achieve similar or even better feature selection effectiveness.

To give a comprehensive performance comparison between 915 PM-CCMB and others, the Friedman test [36] combining with 916 the *post-hoc* tests at a 95% confidence level is used to make 917 a statistical comparison of different algorithms over multiple 918 datasets. The last row of each classifier in Table III shows 919 the average ranks. Note that the traditional feature selection algorithms have better performances against the causal 921 feature selection algorithms except PM-CCMB. It is mainly 922 because all the causal feature selection algorithms are under 923

987

988

990

992

997

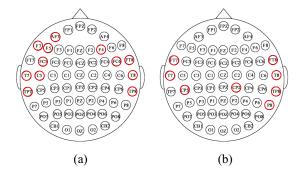


Fig. 8. Profiles of top-20 features selected by PM-CCMB on the (a) Beta and (b) Gamma frequency bands, which is consistent with the previous findings in [35].

924 the faithfulness assumption, which cannot be satisfied in the 925 real-world datasets. However, the proposed cross-check and 926 complement processes can not only avoid the PCMasking 927 phenomenon but also make the PM-CCMB more robust in 928 the unfaithful datasets. Therefore, with the three classifiers, 929 the overall classification accuracy of the proposed PM-CCMB 930 is significantly better than all other algorithms, including the traditional feature selection methods.

Compared with the traditional feature selection methods, 933 PM-CCMB can select useful features and simultaneously 934 explain the causality between features and target. To illus-935 trate the interpretability of PM-CCMB, we collect the features 936 selected by PM-CCMB in all datasets and select the top-937 20 features with the highest frequency, whose positions are 938 illustrated in Fig. 8. As depicted in the figure, the top-939 20 features are all from the Beta and Gamma frequency 940 bands and located at the lateral temporal area except one of 941 them, which is consistent with the previous findings in [35]. 942 These results indicate that PM-CCMB can effectively select 943 the relevant channels containing discriminative information 944 and simultaneously eliminate irrelevant channels for emotion 945 recognition.

VII. CONCLUSION

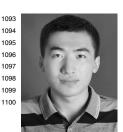
This article introduces the concept of PCMasking to the 948 BN. Based on the PCMasking, we analyze the main reason 949 that causes the lower accuracy in the existing MB discov-950 ery algorithms. To detect the MaskingPCs, we propose the 951 cross-check and complement processes, in which the cross-952 check process can effectively detect the MaskingPCs and the 953 complement process can repair the symmetry (between PC 954 variables) broken by PCMasking phenomenon. On the basis 955 of the cross-check and complement processes, we propose topology-based MB discovery algorithm, CCMB, to find 957 more true variables. To simultaneously improve the time effi-958 ciency, we propose an acceleration methodology for the MB 959 discovery algorithm called PM. Embedding CCMB into PM. 960 we propose PM-CCMB, which can maintain the same accu-961 racy with CCMB and further improve the time efficiency. 962 PM-CCMB is extensively evaluated and compared with the 963 state-of-the-art MB discovery algorithms on different datasets. 964 The results validate that PM-CCMB can shield the influence of PCMasking phenomenon effectively and simultaneously 965 improve the accuracy and time efficiency of MB discovery.

In the future, there are three directions worth further investi- 967 gation. One is to relax the causal sufficiency assumptions [37] 968 and examine the validity of PM-CCMB when applied to real- 969 world feature selection situations. The second direction is to 970 extend PM-CCMB to BN structure learning, especially for 971 large dimensional problems and small samples. The third 972 direction is to develop a joint feature selection and classification algorithm with some Bayesian methods [38], [39].

REFERENCES

- [1] J. Pearl, Probabilistic Reasoning in Intelligent Systems: Networks of 976 Plausible Inference. San Francisco, CA, USA: Morgan Kaufmann, 1998.
- P. Spirtes et al., Causation, Prediction, and Search. Cambridge, MA, 978 USA: MIT Press, 2000.
- [3] C. F. Aliferis, A. Statnikov, I. Tsamardinos, S. Mani, and 980 X. D. Koutsoukos, "Local causal and Markov blanket induction 981 for causal discovery and feature selection for classification part I: 982 Algorithms and empirical evaluation," J. Mach. Learn. Res., vol. 11, no. 1, pp. 171-234, 2010.
- K. Yu, L. Liu, and J. Li, A Unified View of Causal and Non-Causal 985 Feature Selection. [Online]. Available: https://arxiv.org/abs/1802.05844
- [5] D. Margaritis and S. Thrun, "Bayesian network induction via local neighborhoods," in Proc. Adv. Neural Inf. Process. Syst., 2000, pp. 505-511.
- [6] D. Koller and M. Sahami, "Toward optimal feature selection," in *Proc.* 13th Int. Conf. Mach. Learn., Jul. 1996, pp. 284-292.
- I. Tsamardinos, C. F. Aliferis, A. R. Statnikov, and E. Statnikov, "Algorithms for large scale Markov blanket discovery," in *Proc. Florida* 993 Artif. Intell. Res. Soc. Conf., 2003, pp. 376-380.
- [8] I. Tsamardinos, C. F. Aliferis, and A. Statnikov, "Time and sample 995 efficient discovery of Markov blankets and direct causal relations," in 996 Proc. 9th ACM Int. Conf. Knowl. Disc. Data Min., 2003, pp. 673-678.
- C. F. Aliferis, I. Tsamardinos, and A. Statnikov, "HITON: A novel Markov blanket algorithm for optimal variable selection," in *Proc. Amer.* 999 Med. Informat. Assoc. Annu. Symp., 2003, pp. 21-25. 1000
- [10] J. M. Peña, R. Nilsson, J. Björkegren, and J. Tegnér, "Towards scal- 1001 able and data efficient learning of Markov boundaries," Int. J. Approx. 1002 Reason., vol. 45, no. 2, pp. 211-232, 2007.
- S. R. De Morais and A. Aussem, "A novel scalable and data efficient fea- 1004 ture subset selection algorithm," in Proc. Joint Eur. Conf. Mach. Learn. 1005 Knowl. Disc. Databases, 2008, pp. 298-312.
- [12] H. Liu and H. Motoda, Computational Methods of Feature Selection. 1007 Boca Raton, FL, USA: CRC Press, 2007.
- [13] Y. Wang, X. Li, and R. Ruiz, "Weighted general group lasso for gene 1009 selection in cancer classification," IEEE Trans. Cybern., vol. 49, no. 8, 1010 pp. 2860-2873, Aug. 2019. 1011
- J. Cao, Z. Bu, Y. Wang, H. Yang, J. Jiang, and H. Li, "Detecting 1012 prosumer-community groups in smart grids from the multiagent per- 1013 spective," IEEE Trans. Syst., Man, Cybern., Syst., vol. 49, no. 8, 1014 pp. 1652-1664, Aug. 2019.
- [15] Z. Bu, H.-J. Li, C. Zhang, J. Cao, A. Li, and Y. Shi, "Graph 1016 k-means based on leader identification, dynamic game, and opin- 1017 ion dynamics," IEEE Trans. Knowl. Data Eng., to be published. 1018 doi: 10.1109/TKDE.2019.2903712.
- [16] B. Jiang, C. Li, M. D. Rijke, X. Yao, and H. Chen, "Probabilistic feature 1020 selection and classification vector machine," ACM Trans. Knowl. Disc. 1021 Data, vol. 13, no. 2, pp. 1-27, 2019.
- [17] B. Jiang, X. Wu, K. Yu, and H. Chen, "Joint semi-supervised feature 1023 selection and classification through Bayesian approach," in Proc. 33rd 1024 AAAI Conf. Artif. Intell., 2019, pp. 3983-3990.
- [18] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual 1026 information criteria of max-dependency, max-relevance, and min- 1027 redundancy," IEEE Trans. Pattern Anal. Mach. Intell., vol. 27, no. 8, 1028 pp. 1226-1238, Aug. 2005.
- [19] L. Yu and H. Liu, "Feature selection for high-dimensional data: A fast 1030 correlation-based filter solution," in Proc. 20th Int. Conf. Mach. Learn., 1031 2003, pp. 856-863.
- [20] S. Maldonado and R. Weber, "A wrapper method for feature selec- 1033 tion using support vector machines," Inf. Sci., vol. 179, no. 13, 1034 pp. 2208–2217, 2009.

- 1036 [21] Y. Mohsenzadeh, H. Sheikhzadeh, and S. Nazari, "Incremental relevance sample-feature machine: A fast marginal likelihood maximization 1037 1038 approach for joint feature selection and classification," Pattern Recognit., vol. 60, no. 12, pp. 835-848, 2016. 1039
- 1040 [22] B. Xue, M. Zhang, W. N. Browne, and X. Yao, "A survey on evolu-1041 tionary computation approaches to feature selection," IEEE Trans. Evol. Comput., vol. 20, no. 4, pp. 606-626, Aug. 2016. 1042
- Y. Xue, B. Xue, and M. Zhang, "Self-adaptive particle swarm 1043 [23] optimization for large-scale feature selection in classification," ACM 1044 Trans. Knowl. Disc. Data, vol. 13, pp. 1-28, Sep. 2019. 1045
- 1046 [24] I. Guyon et al., "Causal feature selection," in Computational Methods of Feature Selection. Boca Raton, FL, USA: Chapman and Hall, 2007, 1047 1048
- 1049 [25] K. Yu, J. Li, W. Ding, and T. D. Le, "Multi-source causal feature selection," IEEE Trans. Pattern Anal. Mach. Intell., 2019, to be published. 1050 1051 doi: 10.1109/TPAMI.2019.2908373.
- 1052 [26] J.-P. Pellet and A. Elisseeff, "Using Markov blankets for causal structure learning," J. Mach. Learn. Res., vol. 9, no. 7, pp. 1295-1342, 2008. 1053
- 1054 [27] S. Yaramakala and D. Margaritis, "Speculative Markov blanket discovery for optimal feature selection," in Proc. 5th IEEE Int. Conf. Data Min., 1055 1056 Houston, TX, USA, 2005, pp. 809-812.
- S. Fu and M. C. Desmarais, "Fast Markov blanket discovery algorithm 1057 [28] via local learning within single pass," in Proc. Conf. Can. Soc. Comput. 1058 Stud. Intell., 2008, pp. 96-107. 1059
- T. Gao and Q. Ji, "Efficient Markov blanket discovery and its applica-1060 [29] 1061 tion," IEEE Trans. Cybern., vol. 47, no. 5, pp. 1169-1179, May 2017.
- 1062 [30] K. Yu, X. Wu, W. Ding, Y. Mu, and H. Wang, "Markov blanket feature selection using representative sets," IEEE Trans. Neural Netw. Learn. 1063 Syst., vol. 28, no. 11, pp. 2775–2788, Nov. 2017. 1064
- I. A. Beinlich, H. J. Suermondt, R. M. Chavez, and G. F. Cooper, "The 1065 [31] 1066 alarm monitoring system: A case study with two probabilistic inference techniques for belief networks," in Proc. 2nd Eur. Conf. Artif. Intel. 1067 Med., 1989, pp. 247–256. 1068
- 1069 [32] R.-N. Duan, J.-Y. Zhu, and B.-L. Lu, "Differential entropy feature for EEG-based emotion classification," in Proc. IEEE Int. Conf. Neural 1070 1071 Eng., San Diego, CA, USA, 2013, pp. 81–84.
- 1072 [33] W.-L. Zheng and B.-L. Lu, "A multimodal approach to estimating vigilance using EEG and forehead EOG," J. Neural Eng., vol. 14, no. 2, 1073 1074 pp. 17-26, 2017.
- I. Tsamardinos, L. E. Brown, and C. F. Aliferis, "The max-min hill-1075 [34] climbing Bayesian network structure learning algorithm," Mach. Learn., 1076 vol. 65, no. 1, pp. 31-78, 2006. 1077
- W.-L. Zheng, J.-Y. Zhu, and B.-L. Lu, "Identifying stable pat-1078 [35] 1079 terns over time for emotion recognition from EEG," IEEE Trans. Affective Comput., vol. 10, no. 3, pp. 417-429, Jul.-Sep. 2019. 1080 doi: 10.1109/TAFFC.2017.2712143. 1081
- 1082 [36] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," 1083
- J. Mach. Learn. Res., vol. 7, no. 1, pp. 1–30, 2006. K. Yu, L. Liu, J. Li, and H. Chen, "Mining Markov blankets with-1084 [37] out causal sufficiency," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 12, pp. 6333–6347, Dec. 2018. 1085 1086
- 1087 [38] H. Chen, T. Peter, and X. Yao, "Probabilistic classification vector machines," IEEE Trans. Neural Netw., vol. 20, no. 6, pp. 901-914, 1088 Jun. 2009. 1089
- H. Chen, T. Peter, and X. Yao, "Efficient probabilistic classification 1090 [39] vector machine with incremental basis function selection," IEEE Trans. 1091 Neural Netw. Learn. Syst., vol. 25, no. 2, pp. 356-369, Feb. 2014. 1092



Xingyu Wu received the B.Sc. degree from the University of Electronic Science and Technology of China, Chengdu, China, in 2018. He is currently pursuing the M.Sc. degree with the School of Computer Science and Technology, University of Science and Technology of China, Hefei, China.

His current research interests include causal learning, feature selection, and multilabel learning.



Jiang received the B.Sc. degree 1101 from the Chongqing University of Posts and 1102 Telecommunications, Chongqing, China, in 2014, 1103 and the Ph.D. degree in computer science from the 1104 University of Science and Technology of China, 1105 Hefei, China, in 2019.

His current research interests include Bayesian 1107 learning, semisupervised learning, and feature 1108 selection.



Kui Yu received the Ph.D. degree in computer 1110 science from the Hefei University of Technology, 1111 Hefei, China, in 2013.

He is currently a Professor with the Hefei 1113 University of Technology. From 2015 to 2018, he 1114 was a Research Fellow with the University of South 1115 Australia, Adelaide, SA, Australia. From 2013 to 1116 2015, he was a Post-Doctoral Fellow with the School 1117 of Computing Science, Simon Fraser University, 1118 Burnaby, BC, Canada. His current research interests 1119 include causal discovery and machine learning.



Chunyan Miao received the B.S. degree from 1121 Shandong University, Jinan, China, in 1988, and 1122 the M.S. and Ph.D. degrees from Nanyang 1123 Technological University (NTU), Singapore, in 1998 1124 and 2003, respectively.

She is currently a Professor with the School of 1126 Computer Science and Engineering, NTU, and the 1127 Director of the Joint NTU-UBC Research Centre of 1128 Excellence in Active Living for the Elderly (LILY). 1129 Her current research interests include infusing intel- 1130 ligent agents into interactive new media (virtual, 1131

mixed, mobile, and pervasive media) to create novel experiences and dimen- 1132 sions in game design, interactive narrative, and other real-world agent 1133



Huanhuan Chen (M'09-SM'16) received the 1135 B.Sc. degree from the University of Science and 1136 Technology of China (USTC), Hefei, China, in 2004, 1137 and the Ph.D. degree in computer science from the 1138 University of Birmingham, Birmingham, U.K., in 1139 2008.

He is currently a Full Professor with the 1141 School of Computer Science and Technology, 1142 USTC. His current research interests include neu- 1143 ral networks, Bayesian inference, and evolutionary 1144 computation.

Prof. Chen was a recipient of the 2015 International Neural Network 1146 Society Young Investigator Award, the 2012 IEEE Computational Intelligence 1147 Society Outstanding Ph.D. Dissertation Award, the IEEE TRANSACTIONS 1148 ON NEURAL NETWORKS Outstanding Paper Award (bestowed in 2011 1149 and only one paper in 2009), and the 2009 British Computer Society 1150 Distinguished Dissertations Award. He is an Associate Editor of the IEEE 1151 TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS and 1152 the IEEE TRANSACTIONS ON EMERGING TOPICS IN COMPUTATIONAL 1153 INTELLIGENCE.