

---

# ML4C: Seeing Causality Through Latent Vicinity

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Supervised Causal Learning (SCL) aims to learn causal relations from observational  
2 data by accessing previously seen datasets associated with ground truth causal  
3 relations. This paper presents a first attempt at addressing a fundamental question:  
4 *What are the benefits from supervision and how does it benefit?* Starting from seeing  
5 that SCL is not better than random guessing if the learning target is non-identifiable  
6 a priori, we propose a two-phase paradigm for SCL by explicitly considering  
7 structure identifiability. Following this paradigm, we tackle the problem of SCL on  
8 discrete data and propose ML4C. The core of ML4C is a binary classifier with a  
9 novel learning target: it classifies whether an Unshielded Triple (UT) is a v-structure  
10 or not. Starting from an input dataset with the corresponding skeleton provided,  
11 ML4C orients each UT once it is classified as a v-structure. These v-structures are  
12 together used to construct the final output. To address the fundamental question  
13 of SCL, we propose a principled method for ML4C featurization: we exploit the  
14 vicinity of a given UT (i.e., the neighbors of UT in skeleton), and derive features by  
15 considering the conditional dependencies and structural entanglement within the  
16 vicinity. We further prove that ML4C is asymptotically perfect. Last but foremost,  
17 thorough experiments conducted on benchmark datasets demonstrate that ML4C  
18 remarkably outperforms other state-of-the-art algorithms in terms of accuracy,  
19 robustness and transferability. In summary, ML4C shows promising results on  
20 validating the effectiveness of supervision for causal learning.

## 21 1 Introduction

22 The problem of causal learning is to learn causal relations from observational data [13]. The learned  
23 causal relations are typically represented in the form of a Directed Acyclic Graph (DAG), where each  
24 edge in the DAG indicates direct cause-effect relation between the parent node and child node.

25 The methods of causal learning mostly fall into four categories: constraint-based, score-based,  
26 continuous optimization method and functional causal models. Each of these methods takes a given  
27 dataset as input and outputs a DAG but with different criteria. For instance, the DAG should be  
28 consistent with conditional independencies in the data (constraint-based); or it is optimal w.r.t. a  
29 pre-defined score function under either combinatorial constraint (score-based) or continuous equality  
30 constraint (continuous optimization). In a nutshell, these methods can be viewed as *unsupervised*  
31 since they do not access additional datasets associated with ground truth causal relations.

32 A new line of research called *Supervised Causal Learning (SCL)*, on the other hand, aims to learn  
33 causal relations in the supervised fashion: the algorithm has access to datasets associated with  
34 ground truth causal relations, in the hope that such supervision is beneficial to learning causal  
35 relations on newly unseen datasets. Despite several existing works on this direction (see Related  
36 Work), a fundamental question remains unanswered: *How is supervised causal learning possible?*

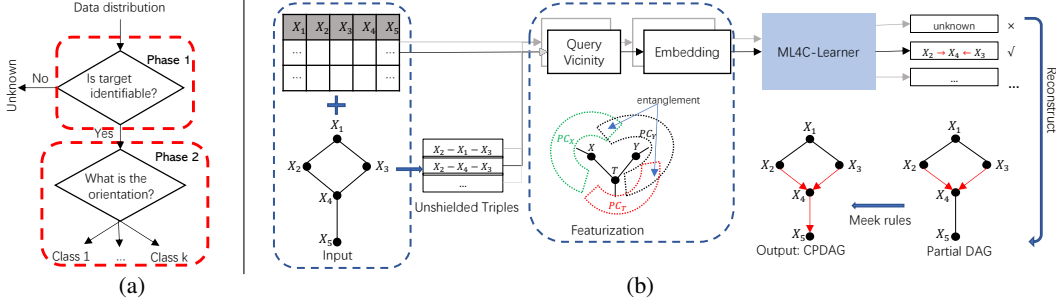


Figure 1: (a) Two-phase paradigm for supervised causal learning. (b) ML4C’s workflow.

37 Specifically, compared with unsupervised causal learning methods, can we gain additional benefits  
 38 from supervision? If the answer is positive, then what are the benefits and how does it benefit?

39 We tackle the problem by first seeing crucial connection between SCL and causal structure identi-  
 40 fiability. Considering the problem of causal learning on discrete data, theorem in [24] states that,  
 41 under standard assumptions (i.e., Markov assumption, faithfulness and causal sufficiency), we can  
 42 only identify a graph up to its Markov equivalence class. Markov equivalence class is the set of  
 43 DAGs having same skeleton and same v-structures, which can be represented by CPDAG (Com-  
 44 pleted Partially Directed Acyclic Graph). Thus, the (un)directed edges in the CPDAG indicate  
 45 (non-)identifiable causal relations. Each non-identifiable edge in CPDAG can be oriented by either  
 46 direction to equivalently fit the observational data. Given an SCL algorithm with learning target as  
 47 the orientation of an edge, we see that it is not better than random guessing (or could be worse due to  
 48 sample bias in training data) to predict any non-identifiable edge since we can assign either  $X \rightarrow Y$   
 49 or  $X \leftarrow Y$  with same input dataset. This statement is applicable to general learning target since an  
 50 SCL algorithm can take different target such as orientation of an edge, the whole DAG, or others.

51 **Proposition 1.** *If the learning target is non-identifiable (i.e., every edge in the target is non-*  
 52 *identifiable) a priori, then SCL is not better than random guessing.*

53 Consequently, we propose and advocate a two-phase paradigm for SCL, as depicted in Figure 1(a):  
 54 phase one corresponds to a binary classification task, where an SCL algorithm needs to classify  
 55 whether a specific learning target is identifiable or not; only if it is classified as identifiable, then we  
 56 go to phase two to classify the specific orientation of the learning target. Following this paradigm,  
 57 we tackle the problem of SCL on discrete data and propose an algorithm ML4C. The core of ML4C  
 58 is a binary classifier with a novel learning target: it classifies whether an Unshielded Triple (UT: a  
 59 triple of variables  $\langle X, T, Y \rangle$  where  $X$  and  $Y$  are adjacent to  $T$  but are not adjacent to each other) is a  
 60 v-structure or not. Starting from an input dataset with the corresponding skeleton provided, ML4C  
 61 orients each UT once it is classified as a v-structure. These v-structures are further used to construct  
 62 a CPDAG as output. Such a single classifier facilitates both learning tasks in the two phases, since an  
 63 identifiable UT implies that it is a v-structure [32] (i.e., up to the partial DAG before applying Meek  
 64 rules [23] which is a standard post processing).

65 To address the fundamental question of SCL, we propose a principled method for ML4C featurization.  
 66 Specifically, we exploit the *vicinity* of a given UT (i.e., the neighbors of UT in skeleton), and derive  
 67 features by considering the conditional *dependencies* and structural *entanglement* within the vicinity.  
 68 We further define discriminative predicate (i.e., a binary predicate function with domain as ML4C’s  
 69 feature set) and prove that there exist weak discriminative predicates and strong discriminative  
 70 predicates (i.e., values of the predicates are one-to-one correspondence with ground truth labels).  
 71 We further prove that ML4C is asymptotically perfect. Last but foremost, thorough experiments on  
 72 benchmark datasets demonstrate that ML4C remarkably outperforms other state-of-the-art algorithms  
 73 w.r.t. accuracy, robustness and transferability. Our main contributions are summarized as follows:

- 74 1. We advocate the two-phase paradigm for SCL with consideration of causal structure identifiability.
- 75 2. We propose an SCL algorithm ML4C, with the following novelties: i) **Learning Target:** The core  
 76 of ML4C is a binary classifier with the orientation of a UT as its learning target to address the two-  
 77 phase tasks simultaneously. ii) **Featurization:** A principled method to exploit vicinity information  
 78 in terms of dependencies and entanglement of a given UT. iii) **Learnability:** We prove that ML4C

79 is asymptotically perfect. iv) **Empirical Performance:** Experiments conducted on benchmark  
80 datasets demonstrate that ML4C remarkably outperforms other state-of-the-art algorithms.

## 81 2 Related Work

82 We divide literature on causal learning into supervised and unsupervised approaches, depending on  
83 whether additional datasets (associated with ground truth causal relations) are accessed (supervised)  
84 or not (unsupervised). In the literature of unsupervised causal learning, constraint-based methods aim  
85 to identify a DAG which is consistent with conditional independencies. The learning procedure of  
86 constraint-based methods first identifies the corresponding skeleton and then conducts orientation  
87 based on v-structure identification [34]. The typical algorithm is PC [31], and there are also PC-  
88 derived algorithms such as Conservative-PC [26], PC-stable [7] and Consistent-PC [19] which  
89 improve the robustness on v-structure identification. Score-based methods aim to find the DAG which  
90 is optimal w.r.t. a pre-defined score function under combinatorial constraint by a specific search  
91 procedure, such as forward-backward search GES [6], hill-climbing [16], integer programming [8], or  
92 by approximate algorithms based on order search [27]. Continuous optimization methods transform  
93 the discrete search procedure into continuous equality constraint: NOTEARS [36] formulates the  
94 acyclic constraint as a continuous equality constraint, it is further extended by DAG-GNN [35] to  
95 support learning non-linear causal relations.

96 SCL emerges from the task of orienting edge in the continuous, non-linear bivariate case under  
97 Functional Causal Model (FCM) formalism. Given a collection of cause-effect samples (dataset  
98  $\sim$  binary label indicating whether  $X \rightarrow Y$  or  $X \leftarrow Y$ ), supervised approaches such as RCC [20],  
99 NCC [21], D2C [4] and Jarfo [12] achieve better performance on predicting pairwise relations (i.e.,  
100 orientation of an edge) than unsupervised approaches such as ANM [14] or IGCI [15]. Differently,  
101 [18] sets the learning target as the whole DAG structure instead of pairwise relation and it is applied on  
102 data which is generated by linear Structural Equation Model (SEM). We summarize the differences in  
103 problem space between ML4C and the other SCL approaches as follows: i) We advocate a two-phase  
104 learning paradigm and emphasize the relationship between identifiability and learnability. Specifically,  
105 presuming additive noise model [14] or linear SEM with non-Gaussian noise [29] provides license  
106 to identifiability thus the aforementioned approaches can be viewed as specific tasks in phase two.  
107 ii) We tackle SCL’s learnability not only via empirical evaluation but also by theoretical analysis to  
108 shed light on the fundamental question of learnability. iii) ML4C deals with discrete data while other  
109 approaches mainly focus on continuous data.

## 110 3 Background

### 111 3.1 Basic Notations

112 A discrete dataset  $D_i$  consists of  $n_i$  records and  $d_i$  categorical columns, which represents  $n_i$  instances  
113 drawn i.i.d. from  $d_i$  discrete variables  $X_1, X_2, \dots, X_{d_i}$  by a joint probability distribution  $P_i$ , which  
114 is entailed by an underlying data generating process, denoted as DAG  $G_i$ .

115 **Markov factorization property:** Given a joint probability distribution  $P$  and a DAG  $G$ ,  $P$  is said to  
116 satisfy Markov factorization property w.r.t.  $G$  if  $P := P(X_1, X_2, \dots, X_d) = \prod_{i=1}^d P(X_i | \text{pa}_i^G)$ ,  
117 where  $\text{pa}_i^G$  is the parent set of  $X_i$  in  $G$ .

118 **Markov assumption:**  $P$  is said to satisfy Markov assumption (or Markovian) w.r.t. a DAG  $G$  if  
119  $X \perp_G Y | Z \Rightarrow X \perp Y | Z$ . Here  $\perp_G$  denotes d-separation, and  $\perp$  denotes statistical independence.  
120 Markov assumption indicates that any d-separation in graph  $G$  implies conditional independence in  
121 distribution  $P$ . Markov assumption is equivalent to Markov factorization property [17].

122 **Faithfulness:** Distribution  $P$  is faithful w.r.t. a DAG  $G$  if  $X \perp Y | Z \Rightarrow X \perp_G Y | Z$ .

123 **Canonical dataset:** We say a discrete dataset  $D$  is canonical if its underlying probability distribution  
124  $P$  is Markovian and faithful w.r.t. some DAG  $G$ .

## 125 3.2 Causal Structure Identifiability

126 Identifiability discusses which parts of the causal structure can in principle be inferred from the  
127 distribution. Below we present the established theory of identifiability on discrete data.

128 **Causal sufficiency:** There are no latent common causes of any of the variables in the graph.

129 **Definition 1** (Markov equivalence). *Two graphs are Markov equivalent if and only if they have*  
130 *same skeleton and same v-structures. A Markov equivalence class can be represented by a CPDAG*  
131 *having both directed and undirected edges. A CPDAG can be derived from a DAG  $G$ , denoted as*  
132  *$CPDAG(G)$ . The theorem of Markov completeness in [24] states that, under causal sufficiency, we*  
133 *can only identify a causal graph up to its Markov equivalence class on canonical data. Therefore, the*  
134 *(non-)identifiable causal relations are the (un)directed edges in the CPDAG. Formally,*

135 **Definition 2** (Identifiability). *Assuming  $P$  is Markovian and faithful w.r.t. a DAG  $G$  and causal*  
136 *sufficiency, then each (un)directed edge in  $CPDAG(G)$  indicates (non-)identifiable causal relation.*

## 137 3.3 ML4C Related Notations

138 **Definition 3** (Skeleton). *A skeleton  $E$  defined over distribution  $P(X_1, X_2, \dots, X_d)$  is an undirected*  
139 *graph such that there is an edge between  $X_i$  and  $X_j$  if and only if  $X_i$  and  $X_j$  are always dependent,*  
140 *i.e.,  $\exists Z \subseteq \{X_1, X_2, \dots, X_d\}$  s.t.  $X_i \perp X_j | Z$ . Skeleton is a statistical concept, which can be*  
141 *obtained prior to facilitating various downstream tasks. Recently, there have been some novel*  
142 *skeleton learning algorithms such as [10]. In particular, skeleton can be used for causal learning:*  
143 *theorem in [32] states that if distribution  $P$  is Markovian and faithful w.r.t. a DAG  $G$ , then skeleton*  
144  *$E$  is the same as the undirected graph of  $G$ .*

145 **Definition 4** (UT). *A triple of variables  $\langle X, T, Y \rangle$  in a skeleton is an unshielded triple, or short*  
146 *for UT, if  $X$  and  $Y$  are adjacent to  $T$  but are not adjacent to each other.  $\langle X, T, Y \rangle$  can be further*  
147 *oriented to become a v-structure  $X \rightarrow T \leftarrow Y$ , in which  $T$  is called the collider.*

148 **Definition 5** (PC). *Denote the set of parents and children of  $X$  in a skeleton as  $PC_X$ , in other words,*  
149  *$PC_X$  are the neighbors of  $X$  in the skeleton. For convenience, if we discuss  $PC_X$  in the context of a*  
150 *UT  $\langle X, T, Y \rangle$ , we intentionally mean the set of parents and children of  $X$  but exclude  $T$ .*

151 **Definition 6** (Vicinity). *We define the vicinity of a UT  $\langle X, T, Y \rangle$  as  $V_{\langle X, T, Y \rangle} := \{X, T, Y\} \cup PC_X \cup$*   
152  *$PC_Y \cup PC_T$ . Vicinity is a generalized version of PC, i.e., the neighbors of  $\{X, T, Y\}$  in the skeleton.*

153 **ML4C’s training set:** The training set is a collection of discrete datasets  $D_1, \dots, D_n$ , where each  
154 dataset  $D_i$  is associated with a ground truth DAG  $G_i$ , such that  $D_i$  is sampled from  $G_i$ .  $G_i$  derives  
155 labels (depends on the chosen learning target), thus  $\{D_i, G_i\}_{i \in \{1, \dots, n\}}$  form ML4C’s training set.  
156 We can sample graphs from DAG space and generate corresponding datasets, thus obtaining training  
157 set in our problem is straightforward.

## 158 4 Approach

### 159 4.1 Overview

160 The core of ML4C is a binary classifier called ML4C-Learner, which takes the orientation of a UT  
161 as its learning target, i.e., it classifies whether an input UT  $\langle X, T, Y \rangle$  is a v-structure (orientation:  
162  $X \rightarrow T \leftarrow Y$ ) or not (orientation remains unknown). Figure 1(b) depicts the overall workflow of  
163 ML4C, which is composed of ML4C-Learner with other important inductive biases. Starting from an  
164 input dataset  $D_i$  with corresponding skeleton  $E_i$ , we first obtain all the UTs from  $E_i$ . Featurization is  
165 then conducted to represent each UT as an embedded vector, which is further fed into ML4C-Learner.  
166 In the inference stage, we obtain all the v-structures which are classified by ML4C-Learner and  
167 reconstruct a partial DAG and then, a CPDAG is output by applying Meek rules on the partial DAG.  
168 In the training stage, the label of each UT is obtained by querying from ground truth DAG  $G_i$ . We  
169 collect labeled data from multiple datasets in ML4C’s training set.

170 **Proposition 2.** *If ML4C-Learner is a perfect classifier, then ML4C outputs correct CPDAG of a*  
171 *canonical dataset (i.e., ML4C is perfect).*

172 By Markov completeness, the set of v-structures is invariant across all Markov equivalent DAGs for  
173 a canonical dataset, and it can fully recover the CPDAG, provided that the skeleton is given. Thus,

174 besides its dedicated role in phase 2, ML4C-Learner also facilitates learning task in phase 1 since an  
 175 identifiable UT implies that it is a v-structure (up to the partial DAG before applying Meek rules).

## 176 4.2 Featurization

177 We propose a principled method for ML4C-Learner’s featurization, which avoids the need of hand-  
 178 crafted features. More importantly, we further prove that ML4C-Learner is asymptotically perfect.

179 **Design Principles:** Our key aspect of featurization is to broaden focus from a specific UT  $\langle X, T, Y \rangle$   
 180 to its *vicinity* and seeking conditional *dependencies* and structural *entanglement* within the vicinity, to  
 181 reveal reliable and robust asymmetry to distinguish v-structure and non-v-structure UTs. Specifically,  
 182 conditional dependencies are the key materials for traditional causal learning methods (e.g., condi-  
 183 tional independences for constraint-based methods), and structural entanglement (e.g.,  $PC_X = PC_T$ )  
 184 are relevant to identifiability: higher entanglement makes the UT less likely to be identifiable.

### 185 • Dependencies within Vicinity

186 **Conditional dependency:** Denoted as  $X \sim Y | \mathbf{Z}$ , which is a non-negative scalar that measures the  
 187 dependence between two random variables  $X$  and  $Y$  given variable set  $\mathbf{Z}$ . Operationally,  $X \sim Y | \mathbf{Z}$   
 188 is composed of two parts, bivariable  $X \sim Y$ , and conditional  $\mathbf{Z}$ . We further extend the definition to  
 189 allow a set of variables in bivariable, and an ensemble (i.e., a set of set) as conditional:

190 **Extended conditional dependency:** Denoted as  $\mathbf{A} \sim \mathbf{B} | \mathcal{Z} := \{X \sim Y | \mathbf{Z} : X \in \mathbf{A}, Y \in \mathbf{B}, \mathbf{Z} \in$   
 191  $\mathcal{Z}\}$ , where  $\mathbf{A}$  and  $\mathbf{B}$  are set of variables, and  $\mathcal{Z}$  is an ensemble. Thus, extended conditional  
 192 dependency is a set of scalars.

193 Within the vicinity of  $\langle X, T, Y \rangle$ , we start from measuring dependencies between  $\{X, PC_X\}$  and  
 194  $\{Y, PC_Y\}$  by conditioning on  $\{T, PC_T\}$ . Intuitively, if  $\langle X, T, Y \rangle$  is a v-structure, conditioning on  
 195  $T$  or  $T$ ’s descendants tends to strengthen the dependency between  $PC_X$  and  $PC_Y$  since the paths  
 196 passing  $X - T - Y$  are unblocked; otherwise, conditioning on  $T$  tends to weaken the dependency  
 197 between  $PC_X$  and  $PC_Y$  because  $T$  blocks the paths passing  $X - T - Y$ . Therefore, such conditional  
 198 dependencies reflect potential asymmetry to distinguish v-structure and non-v-structure. Formally,

199 **Definition 7** (Domain of bivariable). Denoted as  $\mathbb{B} := \{X, PC_X\} \times \{Y, PC_Y\} \equiv$   
 200  $\{X \sim Y, X \sim PC_Y, PC_X \sim Y, PC_X \sim PC_Y\}$ , here symbol  $\times$  is Cartesian product.

201 **Definition 8** (Sepsets). Denoted as  $\mathcal{S} := \{S : X \perp\!\!\!\perp Y | S, S \subset PC_X \cup T, \text{ or } S \subset PC_Y \cup T\}$ . Under  
 202 faithfulness assumption, sepsets  $\mathcal{S}$  is an ensemble where each item is a subset of variables within the  
 203 vicinity that d-separates  $X$  and  $Y$ .

204 **Definition 9** (Domain of conditional). Denoted as  $\mathbb{C} := \{\emptyset, T, \mathcal{P}C_T\} \vee \{\emptyset, S\} \equiv$   
 205  $\{\emptyset, T, \mathcal{P}C_T, S, S \vee T, S \vee \mathcal{P}C_T\}$ , where  $\mathcal{P}C_T := \{I : I \in PC_T\}$  which is an ensemble version  
 206 of  $PC_T$ , and  $S \vee \mathcal{P}C_T := \{S \cup I : S \in \mathcal{S}, I \in \mathcal{P}C_T\}$ . Here symbol  $\vee$  is element-wise union.

207 We exploit the extended conditional dependencies from  $\mathbb{B} \times \mathbb{C}$ , i.e., we pick a bivariable from  $\mathbb{B}$   
 208 and a conditional from  $\mathbb{C}$ , and calculate the extended conditional dependency. There are in total  
 209  $|\mathbb{B}| \times |\mathbb{C}| = 24$  extended conditional dependencies.

210 **Lemma 1.** Sepsets  $\mathcal{S}$  of any UT of a canonical dataset is non-empty. All proofs are available in the  
 211 supplementary material.

212 **Remark 1.** We intend to restrict the sepsets within the vicinity of  $\langle X, T, Y \rangle$ . Lemma 1 shows the  
 213 existence of such d-separation sets within vicinity. Furthermore, searching for all d-separation sets is  
 214 highly time-consuming, thus the computational cost can also be saved drastically.

### 215 • Entanglement within Vicinity

216 Structural entanglement reflects complex structure within the vicinity of  $\langle X, T, Y \rangle$ . Variables  $X, Y$   
 217 and  $T$  can mutually share common neighbors, and their neighbors may also overlap with sepsets  
 218  $\mathcal{S}$ . We call such overlaps structural entanglement. Intuitively, stronger entanglement indicates  
 219 denser structure of vicinity thus making the UT less likely to be identifiable. Therefore, structural  
 220 entanglement is an important aspect for featurization. Specifically, we exploit the overlap coefficient  
 221 [33] to measure the entanglement:

222 **Definition 10** (Overlap coefficient).  $OLP(\mathbf{A}, \mathbf{B}) := |\mathbf{A} \cap \mathbf{B}| / \min(|\mathbf{A}|, |\mathbf{B}|)$ , where  $\mathbf{A}$  and  $\mathbf{B}$  are  
 223 two sets of variables. We extend this formula to support ensemble as input:

224 **(Extended) Overlap coefficient:**  $OLP(\mathbf{A}, \mathcal{S}) := \sum_{i=1}^{|\mathcal{S}|} OLP(\mathbf{A}, S_i) / |\mathcal{S}|$ . Naturally, we consider the  
 225 entanglement in terms of overlap coefficient on each pair of items in domain  $\{PC_X, PC_Y, PC_T, \mathcal{S}\}$ .  
 226 Thus, we use 6 scalars to represent the entanglement within the vicinity of a UT.

227 • **Embedding**

228 We aim to represent the dependencies and entanglement by a feature vector with fixed dimensionality,  
 229 which can be used to train ML4C-Learner. Regarding each extended conditional dependency  $\mathbf{A} \sim$   
 230  $\mathbf{B} | \mathcal{Z} : \mathbf{A} \sim \mathbf{B} \in \mathbb{B}, \mathcal{Z} \in \mathbb{C}$ , it consists of a set of scalars with varied set size across UTs, we  
 231 adopt the kernel mean embedding technique in [30] to represent each  $\mathbf{A} \sim \mathbf{B} | \mathcal{Z}$  as a vector with  
 232 fixed dimensionality. We further modify the embedding algorithm by adding  $\min \{\mathbf{A} \sim \mathbf{B} | \mathcal{Z}\}$  and  
 233  $\max \{\mathbf{A} \sim \mathbf{B} | \mathcal{Z}\}$  as two additional features. We directly use the 6 scalars to represent structural  
 234 entanglement without further transformation. We concatenate all the embedded vectors to form the  
 235 final feature vector, as input for ML4C-Learner.

236 **4.3 Learnability**

237 We have presented ML4C’s featurization and started seeing that conditional dependencies and  
 238 structural entanglement have potential to reveal asymmetry to distinguish v-structure and non-v-  
 239 structure UTs. Now we provide rigorous analysis to show that, for a canonical dataset with sufficient  
 240 samples, ML4C-Learner tends to a perfect classifier. To prove this, we first propose a surrogate object  
 241 called discriminative predicate:

242 **Definition 11** (Discriminative predicate). *A discriminative predicate is a binary predicate function*  
 243 *with domain as ML4C’s feature set. A discriminative predicate can be viewed as a special classifier*  
 244 *with pre-specified form of mechanism (i.e., not learned from data).*

245 **Definition 12** (Weak / Strong discriminative predicate). *Whenever a discriminative predicate takes*  
 246 *the feature vector of a UT as input, a weak discriminative predicate satisfies one of the following*  
 247 *two criteria; a strong discriminative predicate satisfies both: i) it is evaluated to TRUE if the UT is a*  
 248 *v-structure; ii) it is evaluated to FALSE if the UT is not a v-structure.*

249 By definition, a weak discriminative predicate exhibits discriminative power since it is evaluated false  
 250 implies the UT is a non-v-structure (or true implies v-structure). A strong discriminative predicate  
 251 can be viewed as a perfect classifier. Denote  $\{ \mathbf{A} \sim \mathbf{B} | \mathcal{Z} \} > \delta := X \sim Y | \mathcal{Z} > \delta : \forall X \in \mathbf{A}, Y \in$   
 252  $\mathbf{B}, \mathcal{Z} \in \mathcal{Z}$ , then we have:

253 **Lemma 2** (Existence of weak discriminative predicate). *For a canonical dataset with infinite samples,*  
 254 *the following are three weak discriminative predicates: i)  $\{X \sim Y | T\} > 0$ , ii)  $\{X \sim Y | \mathcal{P}C_T\} = 0$*   
 255 *, iii)  $\{PC_X \sim PC_Y | \mathcal{S} \cup T\} > 0$ . Take  $\{X \sim Y | T\} > 0$  as an example,  $\langle X, T, Y \rangle$  is a v-structure  
 256  $\Rightarrow T$  is a collider  $\Rightarrow T$  unblocks  $X$  and  $Y$  through path  $X - T - Y \Rightarrow \{X \sim Y | T\} > 0 \Rightarrow$   
 257  $\min \{X \sim Y | T\} > 0$ , where  $\min \{X \sim Y | T\}$  is a feature of ML4C-Learner since  $X \sim Y \in$   
 258  $\mathbb{B}, \{T\} \in \mathbb{C}$ .*

259 **Lemma 3** (Existence of strong discriminative predicate). *For a canonical dataset with infinite*  
 260 *samples, the following are three strong discriminative predicates: i)  $OLP(T, \mathcal{S}) = 0$ , ii)  $OLP(T, \mathcal{S}) <$*   
 261 *0.5, iii)  $OLP(T, \mathcal{S}) < 1 \wedge \min \{X \sim Y | T \cup \mathcal{S}\} > 0$ .*

262 **CPC/MPC/GLL-MB as special cases of ML4C-Learner:** Predicate  $OLP(T, \mathcal{S}) = 0 \iff \forall S \in$   
 263  $\mathcal{S}, T \notin S$ , which states that the predicate is TRUE if  $T$  is not in any d-separation set of  $X$  and  $Y$ .  
 264 Having correct skeleton provided, this is the criterion of Conservative PC algorithm (CPC) [25] for  
 265 identifying v-structures. Thus, CPC can be viewed as a special case of ML4C by replacing ML4C-  
 266 Learner with such a pre-specified logic;  $OLP(T, \mathcal{S}) < 0.5$  indicates that if more than half of the d-  
 267 separation sets do not contain  $T$ , then the UT is oriented as a v-structure, which is called majority rule  
 268 PC algorithm (MPC) [9]; predicate  $OLP(T, \mathcal{S}) < 1 \wedge \min \{X \sim Y | T \cup \mathcal{S}\} > 0 \Rightarrow \exists S \in \mathcal{S}, T \notin S$   
 269 and  $X$  and  $Y$  are dependent when conditioning on  $T \cup S$ , which is used for GLL-MB [2] to more  
 270 securely identify v-structures. These predicates are with suboptimal performance because only a  
 271 small portion of features are exploited and the overall loss function of training data is disregarded,  
 272 thus in practice when an appropriate machine learning model is adopted, ML4C-Learner achieves  
 273 better performance.

274 **Theorem 1.** *ML4C-Learner tends to a perfect classifier on classifying a canonical dataset with*  
 275 *sufficient samples.*

## 276 5 Evaluation

277 **Benchmark Datasets** We use discrete datasets sampled by all 24 networks from bnlearn repository [28] for evaluation. For each network, we sample 1k, 5k, 10k, 15k, 20k records for use.

279 **ML4C’s Training and Inference** We generate ML4C’s training data synthetically (which is also  
280 used for other SCL competitors). Specifically, 400 unique DAGs are randomly generated by two  
281 models: Erdős-Rényi (ER) model [11] and Scale-Free (SF) model [1], with the number of nodes  
282 ranging from 10 to 1,000. A standard random forward data generation process is applied to obtain  
283 10k observational samples for each graph. We further extract UTs from the 400 DAGs, consisting of  
284 97,010 v-structures (label = 1) and 195,691 non-v-structures (label = 0). We use these instances to  
285 train ML4C-Learner, which is implemented by a XGBoost [5] binary classifier with default hyper-  
286 parameters and we use binary cross-entropy as the loss function. Details on our synthesis procedure,  
287 configurations and implementation of ML4C-Learner are available in the supplementary material.

288 **Competitors** We categorize state-of-the-art causal learning algorithms from two aspects, supervised  
289 vs. unsupervised, and can or cannot take skeleton as input. We choose Jarfo [12], D2C [4], RCC [20],  
290 and NCC [21] as SCL competitors. Same as ML4C, all these algorithms can and do require skeleton  
291 as input. All these algorithms use ML4C’s training set for training but with different learning target  
292 extracted. Regarding unsupervised algorithms, we choose PC [31], Conservative-PC (CPC) [26],  
293 Majority-rule PC (MPC) [7], GLL-MB (GMB) [2], GES [6], Grow-Shrink (GS) [22], Hill-Climbing  
294 (HC) [16], and Conditional Distribution Similarity (CDS) [12]. which can also take skeleton as input.  
295 Lastly, we also compare with DAG-GNN (DGNN) [35], BLIP [27], and GOBNILP (GNIP) [8],  
296 which are unsupervised algorithms but cannot take skeleton as input. All these competitors are  
297 capable of dealing with discrete data. All experiments are done in a Windows Server with 2.8GHz  
298 Intel E5-2680 CPU and 256G RAM. Details are in the supplementary material.

299 **Design** Our evaluation mainly consists of two parts: end-to-end comparison with competitors on  
300 benchmark datasets, and in-depth experiments on ML4C’s learnability. The latter is further divided  
301 into four aspects: i) **Towards a perfect classifier**. As stated in proposition proposition 2, ML4C-  
302 Learner is the core component and we would like to know how far it is from a perfect classifier.  
303 ii) **Reliability** (against weak / strong discriminative predicates). As stated in lemma 2 and 3, there  
304 exist weak and strong discriminative predicates, which have discriminative power and thus are helpful  
305 for ML4C-Learner. Some strong discriminative predicates are equivalent to specific logics of existing  
306 work such as CPC or GLL-MB. Thus, we would like to see how ML4C-Learner takes the advantage  
307 of machine learning, to learn a more reliable classification mechanism (which is also latent and  
308 more sophisticated) than individual weak / strong discriminative predicates. iii) **Robustness** (against  
309 varied sample size). It is known that many causal learning algorithms lack robustness w.r.t sample  
310 noise for finite datasets [20], especially CI tests are error-prone on small samples for constraint-  
311 based algorithms. We would like to evaluate the robustness of ML4C (i.e., the latent classification  
312 mechanism) against varied sample sizes. iv) **Transferability**. It’s important for a machine learning  
313 model to generalize well to various types of testing data which are different from training data, such  
314 as different scale (#nodes), graph sparsity, different generating mechanisms, etc.

315 **Metrics** We use two standard metrics for performance evaluation: Structural Hamming Distance  
316 (SHD) and F1-score. For each dataset, we measure the SHD / F1-score of the output CPDAG (learned  
317 by a specific algorithm) against the ground truth CPDAG. Specifically, SHD is calculated at CPDAG  
318 level, which is the smallest number of edge additions, deletions, direction reversals and type changes  
319 (directed vs. undirected) to convert the output CPDAG to ground truth CPDAG. F1-score is calculated  
320 over identifiable edges. Roughly, F1-score can be viewed as a normalized version of SHD. Now we  
321 present the experiment results:

322 **End-to-End Comparison** Due to page limit, we report SHD and F1-score of all algorithms on  
323 19 large-scale datasets (full results including other 5 smallest and trivial datasets are available in  
324 the supplementary material), as depicted in Table 1. ‘-’ means the algorithm fails on the dataset  
325 (either out-of-memory / exceeds 24 hours execution time / break caused by unknown errors). ML4C  
326 significantly outperforms all other competitors. The average F1-score of ML4C is the highest (0.92,  
327 first column in Table 2). Moreover, ML4C exhibits the most stable performance across all datasets,  
328 its average ranking is  $1.5 \pm 0.7$ , while the second best is GLL-MB (GMB), with average ranking

Table 1: Experiment results for end-to-end comparison with SOTA causal learning algorithms on benchmark datasets. Algorithm names are abbreviated. SHD and F1-score are reported. The last two rows show statistics of rank by SHD and F1-score for all competitors (Note: F1-score is at UT level).

Datasets #nodes/#edges		supervised					unsupervised							no skeleton input			
		ML4C	Jarfo	D2C	RCC	NCC	PC	CPC	MPC	GMB	GES	GS	HC	CDS	DGNN	BLIP	GNIP
child 20/25	SHD F1	<b>0</b> <b>1.0</b>	18 .24	16 .43	18 .33	20 .12	22 .12	13 .00	9 .74	20 .12	15 .47	13 .59	13 .57	18 .34	23 .25	<b>0</b> <b>1.0</b>	<b>0</b> <b>1.0</b>
insurance 27/52	SHD F1	<b>5</b> <b>.89</b>	41 .26	30 .44	34 .42	28 .44	36 .39	34 .00	21 .66	29 .55	34 .46	28 .56	19 .76	36 .36	53 .05	35 .51	14 .82
water 32/66	SHD F1	5 .94	33 .52	43 .34	31 .56	<b>0</b> <b>1.0</b>	4 .97	60 .00	7 .91	8 .87	38 .49	27 .62	38 .46	18 .76	61 .00	65 .20	52 .50
mildew 35/46	SHD F1	6 .87	-	17 .68	25 .50	34 .33	21 .56	-	-	7 .85	<b>3</b> <b>.93</b>	9 .80	23 .64	18 .65	52 .19	36 .41	-
alarm 37/46	SHD F1	<b>1</b> <b>.98</b>	21 .57	26 .44	18 .64	20 .57	20 .57	20 .57	6 .92	17 .64	8 .86	3 .94	21 .66	18 .62	46 .12	17 .82	2 .98
barley 48/84	SHD F1	5 .95	48 .46	55 .38	50 .44	<b>0</b> <b>1.0</b>	3 .96	-	-	8 .91	42 .59	-	34 .72	50 .43	87 .00	60 .48	42 .67
haifinder 56/66	SHD F1	11 .80	47 .21	41 .29	43 .21	<b>0</b> <b>1.0</b>	17 .29	-	-	26 .30	60 .29	-	59 .28	44 .31	76 .00	111 .07	118 .12
hepar2 70/123	SHD F1	<b>0</b> <b>1.0</b>	54 .59	81 .34	59 .54	<b>0</b> <b>1.0</b>	35 .72	27 .81	37 .70	14 .89	46 .75	40 .70	35 .81	75 .39	123 .00	79 .54	61 .68
win95pts 76/112	SHD F1	1 .99	65 .43	51 .54	33 .73	<b>0</b> <b>1.0</b>	8 .95	42 .64	7 .95	5 .97	32 .77	21 .85	16 .91	50 .57	112 .00	103 .47	-
pathfinder 109/195	SHD F1	25 .77	157 .21	145 .29	151 .21	<b>0</b> <b>1.0</b>	150 .29	-	-	147 .30	158 .29	-	168 .28	148 .31	196 .00	241 .07	-
munin1 186/273	SHD F1	<b>10</b> <b>.97</b>	169 .42	154 .47	153 .46	72 .72	86 .71	117 .58	-	84 .72	109 .67	-	233 .26	151 .50	-	257 .42	-
andes 223/338	SHD F1	<b>0</b> <b>1.0</b>	226 .35	209 .41	246 .29	<b>0</b> <b>1.0</b>	4 .99	83 .75	4 .99	5 .98	47 .92	15 .96	38 .92	149 .60	-	175 .76	-
diabetes 413/602	SHD F1	25 .96	220 .62	395 .38	237 .62	48 .96	<b>0</b> <b>1.0</b>	-	-	204 .68	146 .77	-	592 .03	368 .43	-	534 .43	-
pigs 441/592	SHD F1	<b>0</b> <b>1.0</b>	350 .44	332 .46	263 .59	400 .35	400 .35	-	-	268 .56	<b>0</b> <b>1.0</b>	-	532 .18	316 .50	-	6 1.0	-
link 724/1125	SHD F1	<b>0</b> <b>1.0</b>	731 .38	630 .45	638 .45	749 .39	737 .40	-	-	204 .81	324 .80	-	1047 .14	400 .64	-	947 .49	-
munin 1041/1397	SHD F1	72 .95	967 .36	790 .48	816 .44	<b>0</b> <b>1.0</b>	156 .89	-	-	458 .69	661 .62	-	1397 .00	795 .51	-	1599 .29	-
munin2 1003/1244	SHD F1	<b>118</b> <b>.92</b>	554 .60	611 .56	646 .55	1052 .55	898 .30	-	-	536 .57	632 .58	-	1240 .01	753 .49	-	1321 .46	-
munin3 1041/1306	SHD F1	<b>113</b> <b>.92</b>	616 .58	629 .57	688 .54	1048 .25	860 .37	-	-	544 .60	566 .65	-	1306 .00	819 .46	-	1539 .26	-
munin4 1038/1388	SHD F1	<b>126</b> <b>.93</b>	696 .54	658 .56	776 .50	1058 .29	876 .39	-	-	649 .55	618 .64	-	1388 .00	812 .49	-	1627 .28	-
rank(SHD)	mean ±std	<b>1.5</b> <b>0.7</b>	9.1 3.2	8.2 3.7	7.9 2.2	5.1 4.2	6.3 3.7	10.8 2.9	9.5 4.1	4.4 2.4	5.8 2.9	9.6 3.6	8.7 2.7	7.9 2.4	13.3 1.8	10.5 3.5	10.7 4.3
UT-F1	mean ±std	<b>.90</b> <b>.13</b>	.22 .17	.19 .13	.27 .18	.66 .40	.50 .34	.53 .33	.87 .16	.59 .32	.54 .28	.77 .24	.47 .35	.30 .22	.09 .07	.36 .29	.70 .33

Table 2: Reliability: average F1-score of ML4C vs. 8 discriminative predicates extracted from ML4C features on benchmark datasets.

id	ML4C	strong predicates				weak predicates			
		1	2	3	4	1	2	3	4
F1	<b>.92±.20</b>	.77±.31	.52±.27	.38±.25	.66±.27	.72±.25	.61±.29	.73±.30	.55±.27

329 4.4 ± 2.4. Among the competitors, NCC ranks #1 on 8 datasets (note that ML4C ranks #1 on 11  
330 datasets), but its performance fluctuates. Overall it only ranks 5.1 ± 4.2. Last but not least, ML4C  
331 shows high accuracy (F1>0.9) on very large-scale datasets (e.g., medicine datasets ‘munin\*’ [3])  
332 while max(others) ~ 0.6.

333 **Towards a Perfect Classifier** The last row of Table 1 shows the performance of ML4C-Learner  
334 component at UT level by UT-F1 (i.e., F1-score of classifying UTs): such UT level accuracy is  
335 crucial for causal learning on discrete data, since the set of v-structures is invariant across all Markov  
336 equivalent DAGs and it can fully recover the CPDAG. The average F1-score of ML4C-Learner is  
337 0.90 ± 0.13, which shows promising results towards a perfect classifier.



Table 3: Robustness: ML4C is trained on synthetic datasets with sample size = 10k, but tested on benchmark datasets with different sample sizes  $\in \{1k, 5k, 10k, 15k, 20k\}$ .

	size	1k	5k	10k	15k	20k	size	1k	5k	10k	15k	20k	size	1k	5k	10k	15k	20k
SHD	insurance	11	1	5	1	0	water	12	11	5	8	6	mildew	8	5	3	6	1
F1	27/52	.81	.97	.89	.97	1.0	32/66	.86	.87	.94	.89	.93	35/46	.83	.89	.93	.87	.98
SHD	alarm	5	4	0	1	5	barley	13	9	4	8	6	hailfinder	15	15	6	15	13
F1	37/46	.93	.95	1.0	.98	.93	48/84	.88	.93	.97	.92	.94	56/66	.74	.72	.90	.72	.76
SHD	hepar2	8	2	0	1	2	win95pts	7	1	0	1	1	pathfinder	1	7	25	7	1
F1	70/123	.96	.99	1.0	.99	.99	76/112	.96	.99	1.0	.99	.99	109/195	.99	.92	.77	.92	.99
SHD	munin1	32	7	10	9	15	andes	3	2	0	2	0	diabetes	18	28	4	26	27
F1	186/273	.89	.98	.97	.97	.95	223/338	.99	.99	1.0	.99	1.0	413/602	.97	.95	.99	.96	.96
SHD	pigs	0	0	0	0	0	link	88	13	0	0	0	munin	107	76	71	93	87
F1	441/592	1.0	1.0	1.0	1.0	1.0	724/1125	.93	.99	1.0	1.0	1.0	1041/1397	.93	.95	.96	.94	.94
SHD	munin2	117	95	120	110	97	munin3	151	119	113	99	62	munin4	165	130	123	146	133
F1	1003/1244	.92	.93	.92	.93	.93	1041/1306	.90	.92	.92	.94	.96	1038/1388	.90	.92	.93	.91	.93

Table 4: Transferability: ML4C trains/tests both on synthetic datasets with different configurations.

	train	test	SHD	F1	test	SHD	F1	test	SHD	F1	test	SHD	F1
# node	10	10	1.2±2.4	.94±.12	50	4.8±3.4	.95±.03	100	6.6±4.7	.97±.02	1k	50.6±8.4	.97±.00
	50	10	0.4±0.8	.97±.05	50	0.8±1.0	.99±.01	100	4.4±4.7	.98±.02	1k	23.2±5.7	.99±.00
	100	10	0.0±0.0	1.0±.00	50	1.2±1.6	.99±.01	100	4.0±4.6	.98±.02	1k	21.6±4.8	.99±.00
	1k	10	0.4±0.8	.97±.05	50	0.8±1.0	.99±.01	100	1.4±2.3	.99±.01	1k	14.8±8.2	.99±.00
sparsity	1	1	0.8±1.6	.99±.02	2	3.4±2.9	.97±.02	3	3.0±2.5	.98±.01	4	11.4±3.9	.95±.02
	2	1	1.8±1.6	.98±.02	2	2.2±1.7	.98±.01	3	2.2±2.0	.99±.01	4	8.2±2.5	.97±.01
	3	1	1.0±1.3	.98±.02	2	2.2±1.3	.98±.01	3	4.4±3.6	.97±.02	4	4.0±3.2	.98±.01
	4	1	2.4±2.3	.97±.03	2	2.2±1.9	.98±.01	3	3.2±2.7	.98±.02	4	4.8±3.7	.98±.01
sample size	1k	1k	2.8±2.3	.97±.02	5k	2.0±2.2	.98±.02	10k	1.6±2.3	.98±.02	20k	1.0±1.3	.99±.01
	5k	1k	5.2±2.9	.95±.03	5k	1.0±2.0	.99±.02	10k	2.2±3.5	.98±.04	20k	0.6±0.8	.99±.01
	10k	1k	5.2±4.8	.95±.05	5k	1.8±2.7	.98±.02	10k	2.0±3.1	.98±.03	20k	0.6±0.8	.99±.01
	20k	1k	4.8±3.3	.95±.03	5k	2.4±2.6	.98±.02	10k	1.2±1.6	.99±.02	20k	1.0±1.3	.99±.02
gtype	ER	ER	1.0±2.0	.99±.02	SF	2.2±1.6	.98±.01						
	SF	ER	1.6±1.9	.98±.02	SF	2.2±2.4	.98±.02						

338 **Reliability** We manually identify 4 strong discriminative predicates and 4 weak discriminative  
339 predicates and treat each one as a replacement of ML4C-Learner. Table 2 shows the performance of  
340 these predicates. Although most predicates show value on discriminating UTs (e.g., 5/8 predicates  
341 are with  $>0.6$  F1-score), ML4C-Learner has higher performance (average F1-score = 0.92) than each  
342 individual predicate (best average F1-score = 0.77). Thus, it is evident that ML4C-Learner learns a  
343 more reliable classification mechanism, by taking advantage of machine learning techniques.

344 **Robustness** To evaluate robustness, ML4C is trained on synthetic datasets with sample size =  
345 10k, but it is tested on benchmark datasets with different sample sizes: 1k, 5k, 10k, 15k and 20k  
346 respectively. Table 3 shows that ML4C exhibits satisfactory robustness (decrease of F1-score is less  
347 than 0.1) against sample size on most datasets (17/18, except for ‘hailfinder’).

348 **Transferability** To evaluate whether ML4C generalizes well to various types of testing data, we  
349 vary scale (#nodes), graph sparsity, generating mechanism and sample size. ML4C is trained on  
350 a fixed configuration but it is tested with different domains (i.e., data generated under different  
351 configuration). Result is depicted in Table 4, ML4C transfers well on different domains, e.g., even  
352 if it is trained on 10 nodes but tested on 1,000 nodes (last column of the first row in Table 4), the  
353 F1-score only drops 0.02.

## 354 6 Conclusion and Future Work

355 We have proposed a supervised causal learning algorithm ML4C, with theoretical guarantee on  
356 learnability and remarkable empirical performance. More importantly, ML4C shows promising  
357 results on validating the effectiveness of supervision. To make SCL practical in real-world scenarios,  
358 one important direction for future work is to identify reliable and accurate skeleton from data,  
359 considering ML4C requires skeleton as additional input.

## References

- 360  
361 [1] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern*  
362 *physics*, 74(1):47, 2002.
- 363 [2] Constantin F Aliferis, Alexander Statnikov, Ioannis Tsamardinos, Subramani Mani, and Xenofon D  
364 Koutsoukos. Local causal and markov blanket induction for causal discovery and feature selection for  
365 classification part i: algorithms and empirical evaluation. *Journal of Machine Learning Research*, 11(1),  
366 2010.
- 367 [3] Steen Andreassen, Marianne Woldbye, Bjørn Falck, and Stig K Andersen. Munin: A causal probabilistic  
368 network for interpretation of electromyographic findings. In *Proceedings of the 10th international joint*  
369 *conference on Artificial intelligence-Volume 1*, pages 366–372, 1987.
- 370 [4] Gianluca Bontempi and Maxime Flauder. From dependency to causality: a machine learning approach. *J.*  
371 *Mach. Learn. Res.*, 16(1):2437–2457, 2015.
- 372 [5] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd*  
373 *acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- 374 [6] David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine*  
375 *learning research*, 3(Nov):507–554, 2002.
- 376 [7] Diego Colombo and Marloes H Maathuis. Order-independent constraint-based causal structure learning. *J.*  
377 *Mach. Learn. Res.*, 15(1):3741–3782, 2014.
- 378 [8] James Cussens. Bayesian network learning with cutting planes. *arXiv preprint arXiv:1202.3713*, 2012.
- 379 [9] Colombo Diego and H. Maathuis Marloes. Order-independent constraint-based causal structure learning.  
380 *Journal of Machine Learning Research*, 15(116):3921–3962, 2014.
- 381 [10] Rui Ding, Yanzhi Liu, Jingjing Tian, Zhouyu Fu, Shi Han, and Dongmei Zhang. Reliable and efficient  
382 anytime skeleton learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34,  
383 2020.
- 384 [11] P. Erdős and Rényi A. On random graphs. *Publicationes, Mathematicae*, 6:290–297, 1959.
- 385 [12] José AR Fonollosa. Conditional distribution variability measures for causality detection. In *Cause Effect*  
386 *Pairs in Machine Learning*, pages 339–347. Springer, 2019.
- 387 [13] Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical  
388 models. *Frontiers in genetics*, 10:524, 2019.
- 389 [14] Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal  
390 discovery with additive noise models. *Advances in neural information processing systems*, 21:689–696,  
391 2008.
- 392 [15] Dominik Janzing, Joris Mooij, Kun Zhang, Jan Lemeire, Jakob Zscheischler, Povilas Daniušis, Bastian  
393 Steudel, and Bernhard Schölkopf. Information-geometric approach to inferring causal directions. *Artificial*  
394 *Intelligence*, 182:1–31, 2012.
- 395 [16] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press,  
396 2009.
- 397 [17] Steffen L Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.
- 398 [18] Hebi Li, Qi Xiao, and Jin Tian. Supervised whole dag causal discovery. *arXiv preprint arXiv:2006.04697*,  
399 2020.
- 400 [19] Honghao Li, Vincent Cabeli, Nadir Sella, and Hervé Isambert. Constraint-based causal structure learning  
401 with consistent separating sets. In *33rd Conference on Neural Information Processing Systems (NeurIPS*  
402 *2019)*, 2019.
- 403 [20] David Lopez-Paz, Krikamol Muandet, and Benjamin Recht. The randomized causation coefficient. *J.*  
404 *Mach. Learn. Res.*, 16:2901–2907, 2015.
- 405 [21] David Lopez-Paz, Robert Nishihara, Soumith Chintala, Bernhard Scholkopf, and Léon Bottou. Discovering  
406 causal signals in images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
407 *Recognition*, pages 6979–6987, 2017.

- 408 [22] Dimitris Margaritis. Learning bayesian network model structure from data. Technical report, Carnegie-  
409 Mellon Univ Pittsburgh Pa School of Computer Science, 2003.
- 410 [23] Christopher Meek. Causal inference and causal explanation with background knowledge. *arXiv preprint*  
411 *arXiv:1302.4972*, 2013.
- 412 [24] Christopher Meek. Strong completeness and faithfulness in bayesian networks. *arXiv preprint*  
413 *arXiv:1302.4973*, 2013.
- 414 [25] Joseph Ramsey, Jiji Zhang, and Peter Spirtes. Adjacency-faithfulness and conservative causal inference.  
415 *arXiv preprint arXiv:1206.6843*, 2012.
- 416 [26] Joseph Ramsey, Jiji Zhang, and Peter L Spirtes. Adjacency-faithfulness and conservative causal inference.  
417 *arXiv preprint arXiv:1206.6843*, 2012.
- 418 [27] Mauro Scanagatta, Cassio Polpo de Campos, Giorgio Corani, and Marco Zaffalon. Learning bayesian  
419 networks with thousands of variables. In *NIPS*, pages 1864–1872, 2015.
- 420 [28] Marco Scutari. Bayesian network repository. URL <http://www.bnlearn.com>, 2012.
- 421 [29] Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. A linear  
422 non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(10), 2006.
- 423 [30] Alex Smola, Arthur Gretton, Le Song, and Bernhard Schölkopf. A hilbert space embedding for distributions.  
424 In *International Conference on Algorithmic Learning Theory*, pages 13–31. Springer, 2007.
- 425 [31] Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social science*  
426 *computer review*, 9(1):62–72, 1991.
- 427 [32] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and*  
428 *search*. MIT press, 2000.
- 429 [33] M.K. Vijaymeena and K. Kavitha. A survey on similarity measures in text mining. *Machine Learning and*  
430 *Applications: An International Journal*, 3(2):19–28.
- 431 [34] Kui Yu, Jiuyong Li, and Lin Liu. A review on algorithms for constraint-based causal discovery. *arXiv*  
432 *preprint arXiv:1611.03977*, 2016.
- 433 [35] Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks. In  
434 *International Conference on Machine Learning*, pages 7154–7163. PMLR, 2019.
- 435 [36] Xun Zheng, Bryon Aragam, Pradeep Ravikumar, and Eric P Xing. Dags with no tears: Continuous  
436 optimization for structure learning. *arXiv preprint arXiv:1803.01422*, 2018.

## 437 Checklist

- 438 1. For all authors...
- 439 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-  
440 tions and scope? [Yes]
- 441 (b) Did you describe the limitations of your work? [Yes] See §6. ML4C requires skeleton as  
442 additional input, thus we put identifying reliable and accurate skeleton from data as future work.
- 443 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 444 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 445 2. If you are including theoretical results...
- 446 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 447 (b) Did you include complete proofs of all theoretical results? [Yes] Details of proofs are in the  
448 supplementary material due to page limit.
- 449 3. If you ran experiments...
- 450 (a) Did you include the code, data, and instructions needed to reproduce the main experimental  
451 results (either in the supplemental material or as a URL)? [Yes] We have included the main  
452 functionalities of ML4C, synthetic data generator, and ML4C’s training data (include instructions)  
453 for reproducibility. Details are in the supplementary material.
- 454 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?  
455 [Yes] Key information is in content. Details are in the supplementary material.

- 456 (c) Did you report error bars (e.g., with respect to the random seed after running experiments  
457 multiple times)? [Yes] We carefully design experiments, by running experiments multiple times,  
458 including error bars in report to ensure reproducibility. For instance, Table 4 contains error bar  
459 information.
- 460 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,  
461 internal cluster, or cloud provider)? [Yes]
- 462 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 463 (a) If your work uses existing assets, did you cite the creators? [Yes] We compare our algorithm with  
464 15 SOTA competitors and we have cite these work properly. We also use open source benchmark  
465 for evaluation and we also properly cite the creators.
- 466 (b) Did you mention the license of the assets? [N/A]
- 467 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We share  
468 new assets, include main functionalities of our algorithm, and our training data. A private URL  
469 is provided. See details in the supplementary material.
- 470 (d) Did you discuss whether and how consent was obtained from people whose data you're us-  
471 ing/curating? [Yes] Due to page limit, the detailed discussion is available in the supplementary  
472 material.
- 473 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-  
474 tion or offensive content? [N/A] We double-checked that the data we are using/curating contains  
475 no personally identifiable information or offensive content.
- 476 5. If you used crowdsourcing or conducted research with human subjects...
- 477 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?  
478 [N/A]
- 479 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)  
480 approvals, if applicable? [N/A]
- 481 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on  
482 participant compensation? [N/A]