

Bayesian Network Structure Learning Using Causality

Zhen Xu and Sargur N. Srihari

University at Buffalo, The State University of New York, USA
z xu8@buffalo.edu, srihari@cedar.buffalo.edu

Abstract—Bayesian Networks are probabilistic models of data that are useful to answer probabilistic queries. Since they are usually built by hand, algorithms to automatically learn their structure are lacking, particularly for large data sets encountered today. Existing algorithms use either local measures of deviation from independence or global likelihood measures. We tackle this problem from a new perspective using causality, which is a stronger measure than correlation. Integrating both the global and local views, the proposed algorithm learns a high quality Bayesian network without using any score-based searching. Given a partial directed acyclic graph, causal pairs with the highest accuracy are inferred with the fewest number of pairwise causal inferences. Specifically, with discrete data, the χ^2 statistical test is used to identify the most dependent and possible causal pairs. Furthermore, the learned causality is forward-propagated. Experiments on handwriting data show that, besides the ability of causal inference, our algorithm performs better than two previous algorithms, one based on branch-and-bound search, and the other a greedy algorithm using χ^2 tests and a log-loss function. The learned structure not only has lowest loss in representing the data, but also reveals underlying causal relationships which are useful for scientific discovery.

I. INTRODUCTION

Bayesian Networks (BNs) are probabilistic graphical models that represent joint distributions of a set of variables efficiently in a factorized way [1]. A BN is a directed acyclic graph where nodes represent variables and directed edges represent dependency relationships. BNs are widely used in diagnosis, troubleshooting, data mining, etc; some examples are pathologist diagnose lymph node pathologies [1], printer troubleshooting and automobile troubleshooting [2]. Their advantages are in providing a relationship between variables that is easy to understand and in providing a compact form where the amount of data needed to represent the full joint distribution is reasonable.

BNs are largely constructed manually based on expert knowledge. Their use is to infer answers to probabilistic queries. However, with dynamically changing data sets it is useful to automatically learn their structures from data. BNs structure learning is a very active research area, which has three main approaches:

- The constraint-based approach [3], [4] finds a set of conditional independencies and learn the structure of BNs that admit these independencies. Since independencies are only partial properties of the data set, the constraint-based approach usually cannot learn the full structure.
- The score-based approach [1], [5]–[8] searches the

graph space and find the BN structure with the highest score. It is a model selection problem. Because the search space contains a super-exponential number of structures, the search problem is NP-hard [9] and many heuristic algorithms have been proposed.

- The Bayesian model averaging approach generates and averages a set of possible BN structures. Since the number of possible structures can be exponential, some approximation algorithms are devised.

We tackle this problem from a more interesting and fundamental perspective. Roughly speaking, the relationships between two variables are either independent or correlated. The correlation comes either from direct cause and effect relationship between these two variables or from the common cause of these two variables. Instead of using correlation to learn BNs structures the same way as traditional algorithms, we want to use causation to learn the BNs structures. Because causality is more fundamental than correlation, in common sense, the learned structure should be more natural and stable. In addition, it can be served as a good source for scientific discovery. Researchers can search and validate the underlying causality revealed by the learned structure.

II. BN STRUCTURE LEARNING

A BN structure captures independencies that exist in the data to reduce the complexity of the probabilistic model. Methods to construct BNs use a deviation measure from independence and a score function.

A. Deviation Measure

An independence test is used to test independency between two variables given observations. For categorical variables we can use *Pearson's Chi-squared test* to test the null hypothesis of independence. It is a goodness-of-fit test which measures the difference between observed frequency distribution and a theoretical distribution. For two multinomial variables A and B with distributions $P(A = a_i)$ and $P(B = b_j)$, and a data set \mathcal{D} with M samples, we define $O[a_i, b_j]$ is the observed number of samples that $A = a_i$ and $B = b_j$, and $E[a_i, b_j] = MP(a_i)P(b_j)$ is the expected number of samples. The χ^2 test can be defined as:

$$\chi^2(\mathcal{D}) = \sum_{i,j} \frac{(O[a_i, b_j] - E[a_i, b_j])^2}{E[a_i, b_j]}. \quad (1)$$

The bigger test value indicates more dependency between two variables.

B. Score Function

One popular score function for BN learning is the negative log-likelihood, which is known as *log-loss* [1]. Given a data set \mathcal{D} with N i.i.d samples, the log-loss is defined as:

$$l(\mathcal{D}|G) = - \sum_{i=1}^n \sum_{j=1}^N \log P(x_i[j]|\mathbf{x}_{\text{pa}(i)}[j]), \quad (2)$$

where $x_i[j]$ is the value of the i_{th} variable (feature) in the j_{th} sample and $\mathbf{x}_{\text{pa}(i)}[j]$ are values of parents of x_i in the j_{th} sample. It is the loss in terms of using the model to represent the true probability distribution. Since log-probabilities are negative the measure has the convenience of being a positive number.

This score is used by the B & B algorithm and the greedy algorithm described next, but not by the proposed causal algorithm. It is only used to compare the quality of the resulting structure, as described in section VI.

C. Baseline Algorithms

Here we give brief descriptions of two BN structure learning algorithms used to benchmark the performance of the proposed algorithm (Section VI).

1) *Branch and bound (B & B) algorithm:* This algorithm is a score-based approach, which leverages some useful properties of the scoring function to reduce the time and memory costs [10]. Specifically, it uses the minimum description length (MDL) as the score function, which is defined based on log-loss and has the desirable decomposable property. Using this property, it builds cache in $O(n \cdot 2^n)$ time, where n is the number of variables, to avoid redundant computation of log-loss scores. The branch and bound approach makes it an any-time algorithm that can stop at the current best solution while showing an upper bound to global optimum. The time complexity of this algorithm is mainly bound to the cache constructing time. The experiment results are very good for both randomly generated data sets and certain public data sets.

2) *Greedy algorithm:* This uses both the local deviance measure of independence and global log-loss score [11]. It adds and orients one edge at a time while using χ^2 statistics to set the order and log-loss score to decide the orientation. The time complexity of this algorithm is $O(c \cdot 2^q)$, where n is the number of variables and q is the maximum number of parents. Experiments on handwritten data show its good performance. Compared with these two algorithms, our algorithm does not use any global score in the process of structure learning. The log-loss score is only computed with learned BN structure in order to compare the performance among all three algorithms. Because causation is the most important and fundamental factor in correlation, the causal inference is sure to learn a good BN structure, that indirectly implies a good log-loss score. Our experiments in section VI show that the score for learned structure can even be better.

III. CAUSALITY

Broadly speaking, two variables are either independent or correlated. Two variables are independent means the occurrence of one does not influence the probability of the

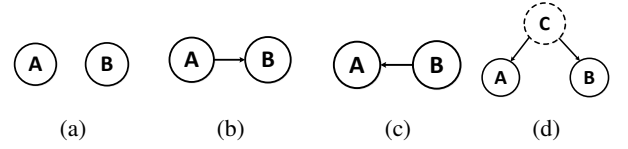


Fig. 1: The possible relationships between events A and B. (a) A is independent with B; (b) A causes B; (c) B causes A; (d) A and B has a common cause, yet they do not cause each other.

other. While two variables are correlated means there is a dependency relationships between these two. Particularly, the causation describe the causal and effect relationship between two variables. There are four possible relationships between two variables A and B [12]: A is independent with B; A causes B; B causes A; A and B are effects of a hidden common cause C, but they do not cause each other [13], as showed in Fig. 1. The independence is illustrated in Fig. 1(a). The correlation includes Fig. 1(b)-(d), while the causation includes Fig. 1(b),(c). from these figures, we can see the causation is more fundamental than correlation and correlation is a broader concept than causation. In most cases, correlation is a good hint to find the causation. While there are many standard statistical independence test to differentiate independency and correlation, it is hard to determine the causation. There are two directions of causal inference from global and local perspectives.

Observe that both causation and correlation are dependent relationships which can be excluded if two variables are independent. PC algorithm [14], which is named by inventors' initials, adopts the constraint-based approach to learn a partial directed acyclic graph (PDAG) which contains directed and undirected edges. This algorithm mainly contains two phases including edge-deleting and edges-orienting. It begins from an undirected complete graph and utilizes the independence test to delete edges. For example, if X and Y are independent, then it deletes edge $X - Y$. If $X|Z$ is conditional independent with $Y|Z$, then it deletes $X - Y$. The edge orienting phase mainly leverages the v-structure and DAG constraints. For example, if there are two edges $X - Y$ and $Y - Z$, X is independent with Z while $X|Y$ is not conditional independent with $Z|Y$, then it orients $X \rightarrow Y \leftarrow Z$. Apply the same procedure until there are no more edges can be oriented by v-structure and DAG constraints, and output the corresponding PDAG, which will still leave lots of edges undirected. The output of this algorithm contains lots of uncertainty and they confined the usefulness of the results for answering probabilistic queries and making predictions. On the other hand, the model-based approaches [15]–[18] infer the causality between two variables by modeling the data generating process. These approaches model the effect variable as a function of the cause variable and the additive noise. In order to break the symmetry and infer the causal relationship between two variables, they made some assumptions such as non-Gaussian noise, non-invertible functional relationship, etc. If the modeled functional relationship is the same as the real data generating process, then the additive noise should be independent with the cause variable. The limitations of model-base approaches come from high

computation complexity. It is only feasible to a very small graph and suffers from multiple hypothesis testing problem.

IV. PROPOSED METHOD

The global and local views of causal inference can be combined to design a computationally efficient BN structure learning algorithm. At the beginning, the constraint-based approach is used to identify the PDAG efficiently, which contains a set of directed and undirected edges. The model-based approach can be used to orient the undirected edges. This is a non-trivial task, because different undirected edges have different causal inference difficulties and accuracies, and they carry different amount causal information which can be used to orient other undirected edges. We need to decide the right order to orient undirected edges and the appropriate way of, With the intuition that the causation is a stronger dependent relationship, we use the χ^2 statistical tests [19] to sort the undirected edges by their dependencies. After identifying the undirected edge with highest dependency, we use the model-based approach to causally orient this edge. with this extra piece of causal information, we apply the causal forward propagation algorithm to other undirected edges which can be oriented further. We repeat this procedure until all edges are causally oriented. Our algorithm uses both the efficiency and stability of the constraint-based method and the causal inference power of the model-based approach to find a BN with small uncertainty. The causal forward propagation algorithm can significant decrease the number of pairwise causal inference and thus save a lot of computations. Moreover, always choosing variable pair with highest dependency improves the accuracy because in common sense the strongest dependency signifies the causality which can be easily identified.

A. Nonlinear Additive Noise Model

The nonlinear additive noise model [16] assumes the nonlinear causal relationship between cause and effect variables. It assumes non-invertible functional relationships between variables to break the symmetry and identify the causality. It models the data generating process as:

$$x_i = f(\mathbf{x}_{\text{Pa}(i)}) + e_i, \quad (3)$$

where f is an arbitrary non-invertible function, $\mathbf{x}_{\text{Pa}(i)}$ is a vector containing all the parents of x_i , and independent noise variable e_i may have arbitrary probability densities $p(e_i)$. If the modeled function is the same as real data generating process, then the noise should be independent with the cause variables.

Given two observed variables x and y . The algorithm first tests whether x and y are statistically independent. If they are not, the algorithm tests weather a model $y = f(x) + e$ is consistent with data. It contains several steps. Firstly, it does a nonlinear regression of y on x to get an estimate \hat{f} of f . Then it calculates the corresponding residuals $\hat{e} = y - \hat{f}(x)$ and test whether \hat{e} is independent of x . If so, it accepts $y = f(x) + e$; if not, it rejects it. The similar test decides whether the reverse model $x = g(y) + e$ fits the data. There are four possible results. Firstly, if x and y are mutually independent, there is no causal relationship between them. Secondly, if they are dependent and both directional models fit the data, then either model may be correct but we do not know which. Thirdly, if

Algorithm 1: CausalBN(V, \mathcal{D})

Input : Vertex Set $V = \{v_1, v_2, \dots, v_n\}$, data set \mathcal{D}
Output : A causal Bayesian network $G^* = \{V, E^*\}$

- 1 Compute PDAG $G_p = \{V, E_p\}$ using Algorithm 2, within E_p there is a list of undirect edges E_u ;
- 2 Compute pairwise χ^2 statistics for edge pairs E_u ;
- 3 Sort E_u in descending order of χ^2 statistics;
- 4 $G^* = G_p$;
- 5 **while** E_u is not empty **do**
- 5 Choose edge $e_t \in E_u$ with the highest χ^2 statistics;
- 6 Infer the causal direct using additive noise model and orient $e_t = \{v_i \rightarrow v_j\}$;
- 7 Apply causal forward propagation, $G^* = \text{CausalForwardProp}(G^*, e_t)$;
- 8 **return** G^* ;

one of the directions is rejected and the other is accepted, a causal relationship is identified. Fourthly, neither model fits the data, which means the real data generating mechanism is more complex and we cannot use the proposed model to describe it.

V. STRUCTURE LEARNING ALGORITHM

The BN structure learning algorithm based on causality has four main steps.

- 1) Use the PC algorithm [14] in Algorithm 2 to identify the PDAG which contains some undirected edges. PC algorithm starts with a complete undirected graph and uses conditional independence to eliminate or orient edges. The edge type for directed pair and undirected pair are expressed as \rightarrow and $-$. Here we define two nodes are adjacent if there is a directed and undirected edge between them.
- 2) Use the χ^2 test to sort the undirected edges based on their dependencies.
- 3) Orient the most dependent undirected pair using nonlinear additive noise model.
- 4) Apply causal forward propagation in Algorithm 3 to this oriented edge to orient other undirected edges recursively.
- 5) Repeat the similar procedure until al edges are oriented.

The complete algorithm is given in Algorithm 1.

The causal forward propagation uses four rules to propagate the causal information when we have a new directed edge. The first rule utilizes directed acyclic property to orient the edges. Note that the circle contains undirected edge does not violate directed acyclic property. The other three rules orient the undirected edge by preventing introducing a new v-structure. Because in PC algorithm, we already use the conditional independency to identify all the possible v-structures completely, so we can use this property to oriented edges. The second rule is straightforward. For the third rule, if we orient $\{v_j - v_k\}$ as $\{v_j \rightarrow v_k\}$, then by DAG constraints, we must orient v_k and v_i as $\{v_i \rightarrow v_k\}$, and $\{v_l \rightarrow v_k\}$ for v_k and v_l . It will produce new v-structure. For the fourth rule, if we orient $\{v_j - v_k\}$ as $\{v_j \rightarrow v_k\}$, then there must be two directed edges

Algorithm 2: Outline of PC-algorithm

Input : Vertex Set $V = \{v_1, v_2, \dots, v_n\}$, significance level α
Output : A PDAG G_P

- 1 Form the complete undirected graph G on the vertex set V ;
- 2 Test conditional and marginal independencies at a given significance level α . Remove one edge between the adjacent pair if there is certain independent relationship between them;
- 3 Orient v-structures by conditional independence information;
while *There are still edges can be oriented by the following rules* **do**
- 4 $\left\{ \begin{array}{l} \text{If } X \rightarrow Y, Y - Z, X \text{ and } Z \text{ are not adjacent, then} \\ \text{oriented } Y - Z \text{ as } Y \rightarrow Z; \\ \text{If } X - Y, \text{ and there is a directed path from } X \text{ to} \\ Y, \text{ then orient } X - Y \text{ as } X \rightarrow Y ; \end{array} \right.$
- 6 **return** G_P ;

$\{v_i \rightarrow v_k\}$ and $\{v_l \rightarrow v_k\}$ to obey the DAG property, it will produce new v-structure.

The time complexity of our algorithm mainly comes from the PC algorithm. In Algorithm 1, PC algorithm in step 1 takes $\frac{n^2(n-1)^{k-1}}{(k-1)!}$ time in the worst case [14], where n is the number of variables, and k is the maximal degree of any vertex. Step 2 takes $O(n^2)$, step 3 takes time $O(n^2 \log(n))$. In addition, the **while** loop runs at most $O(n^2)$ times, within which steps 5 and 6 only take $O(1)$ time. The Algorithm 3 in step 7 needs further consideration. Assume there is no recursion, the first **for** loop in Algorithm 3 runs $O(n^2)$ times, the codes within run $O(1)$ time. The second **for** loop runs $O(n^2)$ times, the codes within run $O(n^2)$ time. The third and fourth **for** loops run $O(n^2)$ times, the codes within run $O(n^4)$ time. So the whole algorithm without recursion runs in $O(n^6)$ time. Note that there are at most $O(n^2)$ undirected edges, so there are at most $O(n^2)$ recursions. Each recursion orients one undirected edge, at the same time, it reduces one **while** loop. So the whole **while** loop takes $O(n^8)$ time. The whole Algorithm 1 takes $\frac{n^2(n-1)^{k-1}}{(k-1)!}$ time in the worst case.

VI. PERFORMANCE EVALUATION

We describe here the performance of the proposed causal structure learning algorithm in terms of log-loss, which measures how well the proposed model fits the data. We use two handwriting data sets, both involving writing of the common word "and":

- 1) Samples from a representative adult population of the United States [20]. There are writing samples of about 1,500 writers containing both hand-printed and cursive writing styles. The data set contains 3,075 cursive "and" samples and 1,135 hand-print "and" samples, both of which have 9 features. The meaning of these features are in [11].
- 2) Handwriting samples of children from the first- to fifth-grades. There are 524 cursive "and" samples and 918 hand-print "and" samples. Each sample has 13

Algorithm 3: CausalForwardProp(G, e)

Input : A PDAG $G = \{V, E\}$ and a newly oriented edge $e = \{v_i \rightarrow v_j\}$
Output : A PDAG G with less uncertainty

- 1 Orient e in E , $E = E - \{v_i - v_j\} + \{v_i \rightarrow v_j\}$;
- 2 Remove e from list E_u ;
- for** *Each edge* $e_1 = \{v_m, v_n\}$ *in list* E_u **do**
 $G_1 = G$;
Orient $e_1 = \{v_m \rightarrow v_n\}$ in G_1 ;
if $\neg isDAG(G_1)$ **then**
3 Orient $e_1 = \{v_n \rightarrow v_m\}$;
4 Apply causal forward propagation, $G =$
CausalForwardProp(G, e_1);
 $G_1 = G$;
Orient $e_1 = \{v_n \rightarrow v_m\}$ in G_1 ;
if $\neg isDAG(G_1)$ **then**
5 Orient $e_1 = \{v_m \rightarrow v_n\}$;
6 Apply causal forward propagation, $G =$
CausalForwardProp(G, e_1);
- for** *Each edge* e_2 *in list* E_p **do**
if *There is an undirected edge* $e_2 = \{v_j - v_k\}$ *and there is no edge between* v_i *and* v_k **then**
7 Orient $e_2 = \{v_j \rightarrow v_k\}$;
8 Apply causal forward propagation, $G =$
CausalForwardProp(G, e_2);
- for** *Each edge* e_3 *in list* E_p **do**
if *There is a undirected edge* $e_3 = \{v_j - v_k\}$, *there are one directed edge* $\{v_l \rightarrow v_j\}$, *and there are edges among* v_k *and* v_i, v_l **then**
9 Orient $e_3 = \{v_k \rightarrow v_j\}$;
10 Apply causal forward propagation, $G =$
CausalForwardProp(G, e_3);
- for** *Each edge* e_4 *in list* E_p **do**
if *There is a undirected edge* $e_4 = \{v_j - v_k\}$, *there are one directed edge* $\{v_l \rightarrow v_i\}$, *there are one edge between* v_l *and* v_k , *and there are no edge between* v_l *and* v_j **then**
11 Orient $e_3 = \{v_k \rightarrow v_j\}$;
12 Apply causal forward propagation, $G =$
CausalForwardProp(G, e_3);
- 13 **return** G ;

features. The meaning of these features can refer to [21].

We compare performance on these two data sets using (i) a state-of-art branch and bound algorithm (B & B) [10], (ii) the greedy algorithm [11] introduced in section II-C, and (iii) the proposed Causal BN learning algorithm.

Table I shows the relative performance of different BN structure learning algorithms using the log-loss measure. A higher score signifies higher probability of the data set given the proposed model. The first two rows give the results with the first data set and the last two with the second data set. The first three columns show the data sets, writing styles, and the number of variables (features). The fourth column corresponds

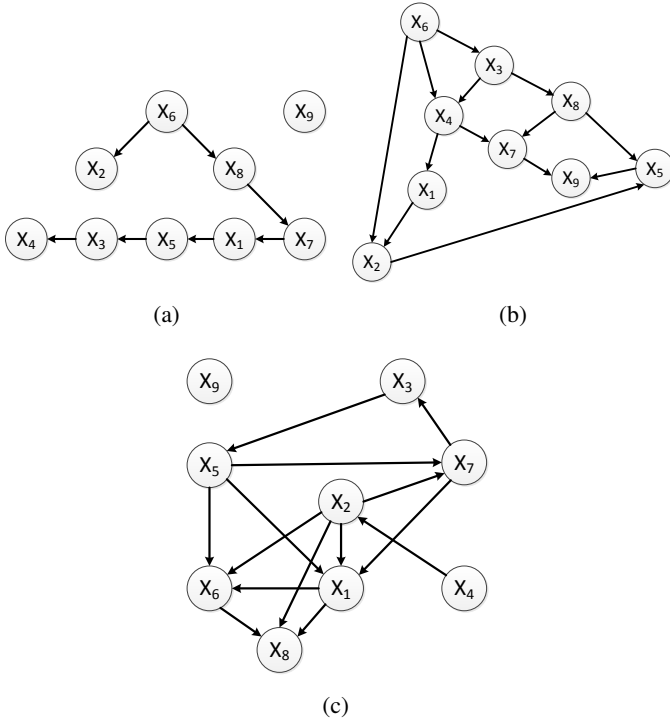


Fig. 2: Bayesian network structures learned from data set 1 (nine variables corresponding to features of "and" that are cursively written): (a) greedy algorithm [11], (b) B & B algorithm, (c) causal algorithm.

TABLE I: Comparison of four BN structure learning methods on four handwriting data sets. Each algorithm is evaluated using the log-loss (Eq. 2) on the data set. Causal algorithm performs best (lowest loss in representing the data).

Data			Algorithm			
Data Set	Type of Data	No. of Vars.	Ind.Vars. (No Edges)	Greedy Algo. [11]	B & B Algo. [10]	Causal Algo.
Set 1	Cursive	9	25994	25329	25642	24988
Set 1	Handprint	9	8059	7898	7301	7094
Set 2	Cursive	12	5316	5142	5139	5008
Set 2	Handprint	12	7825	7004	6976	6956

to log-loss when all variables are independent (no edges in graph). The fifth column contains results from the BN structure learning algorithm in [11]. The sixth column contains results from the B & B algorithm. BNs learned from "and" data set are shown in Fig. 2. From Table I, we can conclude that beside the ability of causal inference, the proposed algorithm performs better in BN structure learning and fits the data sets very well.

Some examples of causality learnt from the data are as follows:

- For cursive 'and' in set 1, the formation of staff of 'd' (X_6), including tented, retracted, looped, etc, decides the formation of terminal stroke of 'd' (X_8), which includes curved up, straight across, curved down, etc.

Formation of "d" staff	tented		retracted		looped		
Formation of "d" terminal	curved up		straight		curved down		no obvious end stroke

Fig. 3: Two variables (X_6 and X_8) and their values in Data Set 1: Cursive 'and'. According to the learnt BN, formation of the staff of d influences the terminal curve.

a-n relationship	a taller		a equal		a smaller	
a-d relationship	a taller		a equal		a smaller	
n-d relationship	n taller		n equal		n smaller	

Fig. 4: Two variables and their values in Data Set 2: Hand-print 'and'. The $a - n$ relationship influences the $n - d$ relationship

The visual examples are showed in Fig. 3. We can see the shape of the character 'd' decides its last stroke.

- For hand-print 'and' in set 2, the 'a'-'d' relationship (x_{11}), including taller, equal, or smaller, decides the 'n'-'d' relationship (x_{12}) with the same values. The visual examples has shown in Fig. 4. Because people usually write a single word in one direction, the 'a'-'d' relationship decides 'n'-'d' relationship.
- For hand-print 'and' in set 2, the 'a'-'n' relationship (x_{10}) decides the 'n'-'d' relationship (x_{12}).

VII. CONCLUSION

A new Bayesian network structure learning algorithm has been proposed. It is deterministic, computationally efficient and uses the principle of causality. Given intuitions and limitations of both global and local causal inference algorithms, we solve a set of challenges to combine them to result in a better algorithm. Specifically, given the causal skeleton, we use the χ^2 statistical test to identify the most dependent pair and perform the causal inference with strong accurate. After that, our causal forward propagation algorithm propagates the newly learned causal information in order to reduce the number of local causal inference. Experiments on handwritten data show that the algorithm provides better statistical models than two previous algorithms, the B & B algorithm and a greedy algorithm. Moreover the resulting structure has the advantage of discovery of causality.

REFERENCES

- [1] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [2] D. Heckerman, J. Breese, and K. Rommelse, "Troubleshooting under uncertainty," *Communications of the ACM*, pp. 121–130, 1994.
- [3] J. Pearl, *Causality: models, reasoning and inference*. Cambridge Univ Press, 2000, vol. 29.
- [4] C. Meek, "Causal inference and causal explanation with background knowledge," in *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 1995, pp. 403–410.
- [5] M. Teyssier and D. Koller, "Ordering-based search: A simple and effective algorithm for learning bayesian networks," *arXiv preprint arXiv:1207.1429*, 2012.

- [6] Y. Barash and N. Friedman, "Context-specific bayesian clustering for gene expression data," in *Proceedings of the fifth annual international conference on Computational biology*. ACM, 2001, pp. 12–21.
- [7] N. Friedman, I. Nachman, and D. Peér, "Learning bayesian network structure from massive datasets: the sparse candidate algorithm," in *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 1999, pp. 206–215.
- [8] G. Elidan, M. Ninio, N. Friedman, and D. Shuurmans, "Data perturbation for escaping local maxima in learning," in *Proceedings of the National Conference on Artificial Intelligence*. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2002, pp. 132–139.
- [9] D. M. Chickering, D. Heckerman, and C. Meek, "Large-sample learning of bayesian networks is np-hard," *The Journal of Machine Learning Research*, vol. 5, pp. 1287–1330, 2004.
- [10] C. P. De Campos, Z. Zeng, and Q. Ji, "Structure learning of bayesian networks using constraints," in *Proceedings of the 26th Annual International Conference on Machine Learning*. ACM, 2009, pp. 113–120.
- [11] M. Puri, S. N. Srihari, and Y. Tang, "Bayesian network structure learning and inference methods for handwriting," in *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*. IEEE, 2013, pp. 1320–1324.
- [12] H. Reichenbach, *The direction of time*. Univ of California Press, 1991, vol. 65.
- [13] F. C. Keil, "Explanation and understanding," *Annual review of psychology*, vol. 57, p. 227, 2006.
- [14] P. Spirtes, C. N. Glymour, and R. Scheines, *Causation Prediction & Search 2e*. MIT press, 2000, vol. 81.
- [15] S. Shimizu, P. O. Hoyer, A. Hyvärinen, and A. Kerminen, "A linear non-gaussian acyclic model for causal discovery," *The Journal of Machine Learning Research*, vol. 7, pp. 2003–2030, 2006.
- [16] P. O. Hoyer, D. Janzing, J. Mooij, J. Peters, and B. Schölkopf, "Nonlinear causal discovery with additive noise models," 2009.
- [17] K. Zhang and A. Hyvärinen, "Distinguishing causes from effects using nonlinear acyclic causal models," in *NIPS 2008 Workshop on Causality*. URL <http://www.cs.helsinki.fi/u/ahyvarin/papers/Zhang09NIPSworkshop.pdf>. Citeseer, 2008.
- [18] —, "On the identifiability of the post-nonlinear causal model," in *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2009, pp. 647–655.
- [19] D. J. Hand, H. Mannila, and P. Smyth, "Principles of data mining (adaptive computation and machine learning)," 2001.
- [20] S. N. Srihari, S.-H. Cha, H. Arora, and S. Lee, "Individuality of handwriting," *Journal of Forensic Sciences*, vol. 47, no. 4, pp. 856–872, 2002.
- [21] M. Puri, S. N. Srihari, and L. Hanson, "Probabilistic modeling of children's handwriting," in *Proc. Document Recognition and Retrieval XXI*, 2014.