

Learning Discrete-valued Bayesian Networks from Mixed Data

First Author · Second Author · Third Author

Received: date / Accepted: date

Abstract Insert your abstract here. Include keywords, PACS and mathematical subject classification numbers as needed.

Keywords Discretization · Bayesian Network · Continuous Variable

1 Introduction

Bayesian networks (Pearl 1988; Koller and Friedman 2009) are an increasingly popular method for modeling uncertainty and causality in science and engineering. They provide an efficient factorization of the joint probability distribution over a set of random variables. Bayesian networks first emerged from artificial intelligence research and have been applied to a wide variety of problems, ranging from decision-making systems (Kochenderfer 2015) to medical diagnoses. In most cases, we assume that all random variables in Bayesian networks are discrete, since many algorithms on Bayesian networks are unable to deal with continuous variables efficiently. However, the assumption that all variables are discrete is often too restrictive. For example, in the decision-making system of autonomous cars, it is imperative to deal with continuous variables such as position and velocity.

There are two methods around this assumption. The first one is to model conditional probability density of each continuous variable by specific families of parametric distributions, then redesign algorithms on Bayesian networks based on these parameters. One successful example is belief propagation in

F. Author
first address
Tel.: +123-45-678910
Fax: +123-45-678910
E-mail: fauthor@example.com

S. Author
second address

Gaussian graphical models (Weiss and Freeman 2001). Nevertheless, for other shapes of particles (Ihler and McAllester 2009) or non-linear functions, these redesigned algorithms are computationally expensive and do not perform well.

The second method around continuous variables is to discretize them. Discretization that learns from data has been developed well and discussed in the fields of machine learning and statistics for many years (Dougherty et al. 1995; Kerber 1992; Holte 1993; Fayyad and Irani 1993). Most of these discretization methods are designed for classification problems. They search the best discretization policy of a continuous attribute by considering its interaction with the class variable of interest. However, these discretization methods do not apply to continuous variables in Bayesian networks. In Bayesian networks, interactions and dependencies between variables are determined by graph structure. Therefore, if a method discretizes continuous variables according to graph structure, instead of assigning one variable as a class variable, then it would be more appropriate. There is some research on discretizing continuous variables in naive Bayesian networks and tree augmented networks (Friedman et al. 1997), but only a few discretization methods on general Bayesian network have been proposed (Friedman and Goldszmidt 1996; Kozlov and Koller 1997; Monti and Cooper 1998; Stech and Jaakkola 2007).

The discretization technique proposed by Friedman and Goldszmidt (1996) is the most well-known one among these few methods. It is based on minimum description length principle (MDL): optimal discretization policy minimizes the description length of Bayesian network and data information. If there is only one continuous variable in Bayesian network and other variables are discrete, this MDL discretization method takes running time $O(N^3 + n_c v_{max} (n_c^p)_{max} \cdot N^2 + v_{max}^{n_p} \cdot N^2)$, where N is the number of data instances for learning the discretization, n_p and n_c are the numbers of parent and children variables for the continuous variable, v_{max} is the largest cardinality number over all variables in Markov blanket, and $(n_c^p)_{max}$ is the largest number of parent variables of the continuous variable's children. If there are multiple continuous variables in the network, the method iterates over all continuous variables until convergence. In our test on real world data, the iterations converge in a few cycles. However, the number of intervals after discretization is too small.

In this paper, we propose a new discretization method for continuous variables in Bayesian networks by learning from mixed data. This method is a generalized version of Boullé (2006) and Lustgarten et al. (2011), which are both discretization methods of a continuous attribute with one class variable. We begin our method with the assumption that only one variable in a given Bayesian network is continuous and other variables are discrete. Furthermore we assume that the network structure is known in advance. Under this situation, we look for the most probable discretization policy M given the data D on other discrete variables. That is to say, the desired policy is $\arg_M P(M | D)$. With Bayes' rule, $P(M | D)$ can be rewritten as $P(M) \cdot P(D | M) / P(D)$, which is proportional to $P(M) \cdot P(D | M)$. Usually, $P(D | M)$ increases as the number of discretized intervals increases, since more intervals can provide

more accuracy. $P(M)$, on the other hand, is designed to decrease as the number of intervals rises. As a result, we can automatically determine the number of intervals after discretization by maximizing $P(M) \cdot P(D | M)$. In addition, the $P(D | M)$ can be factorized according to Bayesian network structure. With the proposed priors, we are able to restrict the number of discretized intervals so that it will not exceed the largest cardinality number of variables in the Markov blanket by too much. This is important, since for most algorithms on Bayesian networks, their running times exponentially depend on the cardinality of variables. Another advantage of our method is the running time for learning. If there is only one continuous variable, the running time is $O(n_c v_{max} \cdot N^2 + v_{max}^{n_p} \cdot N^2)$, which is significantly shorter than the running time of MDL principle discretization (Friedman and Goldszmidt 1996). With the approximation proposed in this paper, we can further reduce running time to $O(v_{max}(n_p + n_c)N^2)$.

For a Bayesian network with multiple continuous variables, we apply the discretization method discussed above iteratively. In the beginning, we prediscritize each continuous variable by equal-width method. The number of intervals for equal-width discretization is equal to largest cardinality of discrete variables. Then we iterate the one-variable discretization on each continuous variable in the following order: from the variable with highest topological order (leaves) to the variable with lowest topological order (root). Each time we finish a discretization procedure on one variable, we store the discretization result for later iterations. Experiments on real data show that with the same iteration order, our discretization method provides better discretization results than MDL principle method in terms of likelihood. The MDL principle method is easily stuck at local minima and also discretizes variables into too few intervals.

Finally, we can combine our new discretization method with K2 structure learning algorithm (Cooper and Herskovits 1992). We first prediscritize all continuous variables before K2 learning, since the K2 algorithm requires all variables to be discrete. Each time an edge is added in the K2 algorithm procedure, we rediscritize the relevant variables in the learned Bayesian network and store the result for next K2 algorithm iteration. By this principle, we are able to learn a discrete-valued Bayesian network from mixed data.

The rest of the paper is organized as follows: related works, including MDL principle discretization (Friedman and Goldszmidt 1996) and MODL discretization by Boullé (2006), are summarized in Section 2. We prepare the preliminaries and notations in Section 3. In Section 4, we introduce the new discretization method for the case that only one variable is continuous, including the derivation of objective function and the algorithms. Section 5 is a discussion about multiple continuous variables in a Bayesian network and how to discretize them. The combination of our proposed discretization and K2 structure learning is also introduced. Finally, in Section 6, we apply the proposed discretization method to real-world data and show the result.

2 Related Work

In this section we review two related works: MDL principle discretization (Friedman and Goldszmidt 1996) and MODL discretization (Boullé 2006). The former is the most famous discretization method for continuous variables in Bayesian networks, and we will compare it with our proposed method in Section 6. The latter is a discretization method for one continuous attribute according to a target class. Our proposed method is a generalization of this method. Note that MODL discretization is a Bayesian approach. The asymptotical equivalence between MDL approach and Bayesian approach has been examined in (Vitényi and Li 2000).

2.1 MDL Principle Discretization

The MDL principle was first proposed by Rissanen (1978). It provides a way to compress data, that is to say, describe data with fewer number of symbols than the number of symbols needed to describe the data literally (Grünwald 2007). MDL chooses a model that trades off goodness-of-fit for the complexity of model. Therefore it has an advantage that automatically avoids overfitting. In the case of Bayesian networks, (Friedman and Goldszmidt 1996) applied the MDL concept to determine the number of intervals of each continuous variable after discretization and also the positions of discretization boundaries. Here is the mechanism: MDL principle discretization selects a discretization policy that minimizes sum of description lengths of discretized Bayesian network and the necessary information that recovers original data from discretized data. If there is only one continuous variable and other variables are discrete, the objective function is

$$\begin{aligned} & \frac{1}{2} \log(N) \{ \|\Pi_{X_i}\| (\|X_i^*\| - 1) + \sum_{j, X_i \in \Pi_{X_i}} \|\Pi_{X_j^*}\| (\|X_j\| - 1) \} + \\ & \log(\|X_i^*\|) + \log\left(\frac{N_i - 1}{\|X_i^*\| - 1}\right) - N \cdot [I(X_i^*, \Pi_{X_i}) + \sum_{j, X_i \in \Pi_{X_i}} I(X_j, \Pi_{X_j^*})], \end{aligned} \quad (1)$$

where function $I(\mathbf{A}, \mathbf{B}) = \sum_{\mathbf{a}, \mathbf{b}} \hat{P}_D(\mathbf{a}, \mathbf{b}) \log \frac{\hat{P}_D(\mathbf{a}, \mathbf{b})}{\hat{P}_D(\mathbf{a}) \hat{P}_D(\mathbf{b})}$ is the mutual information between two sets of variables \mathbf{A} and \mathbf{B} , $\|\mathbf{X}\|$ means the cardinality of variable \mathbf{X} and X_i^* is the discretized version of continuous variable X_i . The optimal discretization policy can be found by dynamical programming. If there are multiple continuous variables in Bayesian network, we apply the algorithm on one variable at a time and iterate over all continuous variables. While iterating, only one variable is treated as continuous and other continuous variables are discretized based on initial prediscrretization or the discretization result of previous iteration.

For the above discretization processes, the network structure is known in advance. If the network structure is not given, we can alternate between structure learning and discretization learning. We start with some discretization,

and learn a network structure given this discretization. Then, we rediscretize based on the learned network. This cycle continues until no improvement in objective function is made.

2.2 MODL Discretization

MODL discretization Boullé (2006) is a Bayesian-approach discretization. It discretizes a continuous feature according to a class variable. The best discretization model is found by maximizing the probability $P(\text{Model} \mid \text{Data})$ of the model given the data. Using the Bayesian rule and since $P(\text{Data})$ is a constant while varying the discretization model, this is equivalent to maximize:

$$P(\text{Model}) \cdot P(\text{Data} \mid \text{Model}). \quad (2)$$

Once the prior distribution of the discretization model is fixed, each term in Eq (3) can be evaluated. Usually, $P(\text{Data} \mid \text{Model})$ increases as the number of intervals of discretization model increases, since more intervals means that the model contains more information and thus it is easier to reproduce the original data. On the other hand, $P(\text{Model})$ decreases as the number of intervals increases. Thus, maximizing $P(\text{Model}) \cdot P(\text{Data} \mid \text{Model})$ provides a trade-off to determine number of intervals after discretization.

The MODL method adapts dynamical programming to find the optimal discretization model, and it takes running time $O(N^3 + \|\text{Class}\| \cdot N^2)$, where N is number of data instances. Lustgarten et al. (2011) provide a different prior to evaluate $P(\text{Model})$ term. This prior could include the assumption of uniform prior probabilities over discretization in MODL as a special case. Additionally, it reduces runing time to $O(\|\text{Class}\| \cdot N^2)$.

3 Preliminaries

In this section we provide a brief review of Bayesian networks, including the factorization of joint probability distribution, sampling from a given network, and structure learning. These concepts will be used in later sections.

3.1 Bayesian Network and Structure Learning

A Bayesian network B is defined by a pair (G, Θ) , where $G = (X, E)$ is a directed acyclic graph whose nodes correspond to a set of random variable $X = \{X_1, X_2, \dots, X_n\}$, and whose edges E represents probabilistic dependencies among nodes. The graph structure G encodes the Markov property: each node X_i is independent of its non-descendants given its parents in G . The second component of B , namely Θ , contains a set of parameters that qualify the network. Elements of Θ take the form $\theta_{x_i \mid \Pi_{x_i}} = P(x_i \mid \Pi_{x_i})$ for each possible value x_i of X_i , and Π_{x_i} of Π_{X_i} (the set of parents of X_i in G). Applying

Markov property, we can represent the multivariate joint distribution over X as

$$P_B(X_1, \dots, X_n) = \prod_{i=1}^n P_B(X_i | \Pi_{X_i}) = \prod_{i=1}^n \theta_{x_i | \Pi_{x_i}} \quad (3)$$

For a given Bayesian network, we can generate its data instances by forward sampling, i.e., sampling variables one by one in a topological order (see CORMEN et al. 2009, chap. 22). Given an instance of X_i 's parent set, Π_{X_i} , values of X_i can be sampled according to the conditional probability table $P(X_i | \Pi_{X_i})$. Furthermore, this parent-child sampling order can be reversed if we know the marginal probability of child variable $P(X_i)$. By Bayes rule,

$$P(X_i | \Pi_{X_i}) \cdot \prod_{j=1}^{Pa(X_i)} P(\{\Pi_{X_i}\}_j) = P(X_i) \cdot P(\Pi_{X_i} | X_i), \quad (4)$$

where $Pa(X_i)$ is the number of parents of X_i , and $\{\Pi_{X_i}\}_j$ is the j th parent of X_i . We can first sample X_i by $P(X_i)$, then sample all the parent variables simultaneously by $P(\Pi_{X_i} | X_i)$. That is to say, X_i becomes the starting point for sampling.

If we don't know the structure of the Bayesian network in advance, then we need to learn it from data. Roughly speaking, there are three approaches to learn a Bayesian network structure from data (see Koller and Friedman 2009, chap. 18): constraint-based structure learning, score-based structure learning, and Bayesian model averaging. Here we review K2 structure learning algorithm (Cooper and Herskovits 1992), which is one of the most successful score-based structure learning methods. Similar to most structure learning algorithms, the K2 algorithm requires all variables to be discrete. The score of a learned network is defined as $\prod_i f(X_i, \Pi_{X_i})$, where

$$f(X_i, \Pi_{X_i}) = \prod_{j=1}^{|\Pi_{X_i}|} \frac{(|X_i| - 1)!}{(N_{ij} + |X_i| - 1)!} \prod_{k=1}^{|X_i|} \alpha_{ijk}! \quad (5)$$

$|X_i|$ is the cardinality of variable X_i . $|\Pi_{X_i}|$ is number of all possible instantiations of the parent variables of X_i , i.e., $|\Pi_{X_i}| = \prod_{Y \in \Pi_{X_i}} |Y|$. α_{ijk} is the number of instances in a data set that variable X_i is instantiated with its k th value, and the parent of X_i are instantiated with the j th value of Π_{X_i} . $N_{ij} = \sum_{k=1}^{|X_i|} \alpha_{ijk}$.

K2 algorithm searches the network structure with the highest score, which can be interpreted as the most probable network with the given data. In addition, K2 requires a topological order of variables to be known before scoring can start. This constraint can prevent cycles from being introduced. The searching for high-score network is an iterative process. There is no way to find the optimal network directly, since there are $2^{O(n^2)}$ possible structures, where n is number of variables. In order to compensate for this, we run K2 algorithm many times, and each time we start with a different order of variables. The network with the highest score in iterations is the desired one.

3.2 Discretization Policy

A Discretization Policy $M_C = \langle t_1, t_2, \dots, t_{k-1} \rangle$ on a continuous variable C is a mapping from \mathbf{R} to $\{1, 2, 3, \dots, k\}$ such that

$$M_C(x) = \begin{cases} 1, & \text{if } x < t_1. \\ i, & \text{if } t_{i-1} \leq x < t_i. \\ k, & \text{if } t_k \leq x \end{cases} \quad (6)$$

That is to say, the policy discretizes the continuous variable C into k intervals. Furthermore, we assume that the discretization edge t_i can only take values on middle points of two consecutive values of data of C . By this assumption, we can obtain an integer representation of M_C as follows: we sort the data of continuous variable C in an increasing order, $\{c_1, c_2, \dots, c_N\}$, and if $t_i = (c_{s_i} + c_{s_i+1})/2$ for $i = 1, 2, \dots, k$, then $M_C = \langle t_1, t_2, \dots, t_{k-1} \rangle \equiv [n_1, n_2, \dots, n_k]$, where $n_1 = s_1$ and $n_i = s_i - s_{i-1}$ for $i = 2, \dots, k$.

4 Discretize One Continuous Variable

In this section, we consider the case that only one variable in a Bayesian network is continuous and the others are discrete. Before we formulate the objective function, we introduce some notations that will make the following calculation easier.

4.1 Notations

Assume the continuous variable is X . It has n_p parent variables, $\Pi_X = \{P_1, P_2, P_3, \dots, P_{n_p}\}$. It also has n_c child variables $\{C_1, C_2, \dots, C_{n_c}\}$. Each child variable C_i of X has other parent variables, which forms a set \mathbf{S}_i . Before we do the discretization, we sort the data instances by its attribute of X .

Let $D \cup D_X$ be a dataset of N instances which we plan to learn discretization from. D_X only contains the data of X . D contains the data of all variables except X . Assume $D \cup D_X$ has been sorted in the ascending order of D_X . In the following content, $D_{\mathbf{Y}}$ means the data instances of a set of variable \mathbf{Y} . Let $D_X = \{x_1, x_2, x_3, \dots, x_N\}$. In order to make the discretization mechanism more understandable, we temporarily assume there is no repeated values in D_X , i.e., $x_1 < x_2 < \dots < x_N$. We will remove this assumption later in the algorithm section. Therefore, a discretization policy M on X can be written as $M = [n_1, n_2, \dots, n_k]$, where k, n_1, n_2, \dots , and n_k are positive integers satisfying $N = \sum_{i=1}^k n_i$.

4.2 Priors and Objective Function

We have the following four priors for the discretization model M that enable us to evaluate $P(M)$ and $P(D | M)$:

1. For a discretization edge locating at $(x_i + x_{i+1})/2$ has probability

$$1 - \exp(-M \cdot \frac{x_{i+1} - x_i}{x_N - x_1}), \quad (7)$$

where M is the largest cardinality number of discrete variables in X 's Markov blanket.

2. For a given interval of M , every distribution of Π_X is equiprobable.
3. For each pair of (C_i, \mathbf{S}_i) and a given interval of M , every distribution of C_i with a value of \mathbf{S}_i is equiprobable.
4. The distributions of each C_i and Π_X in each interval of M are independent from each other.

Owing to these priors, we are able to evaluate $P(M)$ and $P(D | M)$. For the former, we have

$$P(M) = \prod_{i=1}^k \{ [1 - \exp(-M \cdot \frac{x_{s_i+1} - c_{x_i}}{x_N - x_1})] \cdot \exp(M \cdot \frac{x_{s_i} - x_{s_{i-1}+1}}{x_N - x_1}) \}, \quad (8)$$

where $s_i = \sum_{j=1}^k n_j$ and $s_k = N$, $s_0 = 0$. Before evaluating $P(D | M)$, notice that

$$P(D | M) \propto P(D_{\Pi_x} | M) \cdot \prod_{i=1}^{n_c} P(D_{C_i} | M, D_{\mathbf{S}_i}). \quad (9)$$

The graph structure allows us to only consider the interactions between X and other variables in its Markov blanket and also make the factorization. For example, in Figure 1., the corresponding $P(D | M)$ is factorized as

$$P(D_{P_1, P_2, P_3} | M) \cdot P(D_{C_1} | M, D_{\mathbf{S}_1}) \cdot P(D_{C_2} | M, D_{\mathbf{S}_2}). \quad (10)$$

The concept behind the factorization is similar to sampling in a Bayesian network. Once we have the distribution of X and sample X for it, then we can further sample Π_x . Since the samples of variables in Π_x are not independent, we can not further factorize $P(D_{\Pi_x} | M)$ into $\prod_i P(D_{P_i} | M)$. Similarly, once we have a set of samples of X and \mathbf{S}_i , we can sample C_i . Since the sampling processes of child variables are independent from each other, we have the product $\prod_{i=1}^{n_c} P(D_{C_i} | M, D_{\mathbf{S}_i})$.

For the following part, we evaluate $P(D_{\Pi_x} | M)$ and $P(D_{C_i} | M, D_{\mathbf{S}_i})$ based on discretization policy $M = [n_1, n_2, \dots, n_k]$ and dataset D .

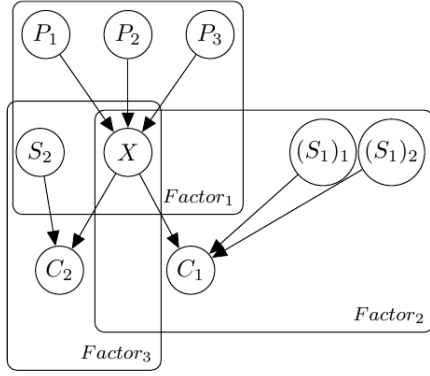


Fig. 1 Factorization of $P(D | M)$

4.2.1 Evaluate $P(D_{\Pi_X} | M)$

Assume that $J_P = \prod_{i=1}^{n_p} \|P_i\|$. Then we have:

$$P(D_{\Pi_X} | M) = \prod_{i=1}^k \frac{1}{\binom{n_i + J_P - 1}{J_P - 1}} \frac{1}{\frac{n_i!}{n_{i,1}^{(p)}! n_{i,2}^{(p)}! \dots n_{i,J_P}^{(p)}!}}, \quad (11)$$

where $n_{i,j}^{(p)}$ is the number of instances in i th discretized interval with j th value of Π_X . Note that $n_i = \sum_{j=1}^{\|\Pi_X\|} n_{i,j}^{(p)}$. The two factors on RHS comes from the second prior: all distributions of values of Π_X in a given interval are equiprobable. According to the forth prior, the distribution in each interval is independent, so we multiply all the two factors together.

4.2.2 Evaluate $P(D_{C_i} | M, D_{S_i})$

Assume for each pair of (C_j, S_j) , we have $\|C_j\| = J_j$ and $\|S_j\| = L_j = \prod_{v \in S_j} \|v\|$. Therefore,

$$P(D_{C_j} | M, D_{S_j}) = \prod_{i=1}^k \prod_{l=1}^{L_j} \frac{1}{\binom{n_{i,l} + J_j - 1}{J_j - 1}} \frac{1}{\frac{n_{i,l}!}{n_{i,1,l}^{(j)}! n_{i,2,l}^{(j)}! \dots n_{i,J_j,l}^{(j)}!}}, \quad (12)$$

where $n_{i,m,l}^{(j)}$ is the number of instances in i th interval with m th value of C_j and l th value of S_j , $n_{i,l} = \sum_{m=1}^{J_j} n_{i,m,l}^{(j)}$, and $n_i = \sum_{l=1}^{L_j} n_{i,l}$. The two factors on RHS comes from the third prior: all distribution of values of C_j in a given interval and with a given value of S_j are equiprobable. According to the forth prior, these distributions are independent from each other, therefore we

multiply all factors. If $S_j = \emptyset$, then Equation 13. is equivalent to

$$P(D_{C_j} | M) = \prod_{i=1}^k \frac{1}{\binom{n_i + J_j - 1}{J_j - 1}} \frac{1}{\frac{n_i!}{n_{i,1,\emptyset}^{(j)}! n_{i,2,\emptyset}^{(j)}! \dots n_{i,J_j,\emptyset}^{(j)}!}}, \quad (13)$$

where $n_{i,m,\emptyset}^{(j)}$ is the number of instances in i th interval with m th value of C_j , and $\sum_{m=1}^{J_j} n_{i,m,\emptyset}^{(j)} = n_i$.

With Equation 8,9,11,12, and 13, now we are able to write down the objective function. Instead of maximizing $P(M) \cdot P(D|M)$, we minimize its log-inverse for convenience. The objective function is

$$\begin{aligned} & \sum_{i=1}^{k-1} -\log(1 - \exp(-M \cdot \frac{x_{s_{i+1}} - x_{s_i}}{x_N - x_1})) + \sum_{i=1}^k M \cdot \frac{x_{s_i} - x_{s_{i-1}+1}}{x_N - x_1} + \\ & \sum_{j=1}^{n_c} \sum_{i=1}^k \sum_{l=1}^{L_j} \left[\log \binom{n_{i,l} + J_j - 1}{J_j - 1} + \log \left(\frac{n_{i,l}!}{n_{i,1,l}^{(j)}! n_{i,2,l}^{(j)}! \dots n_{i,J_j,l}^{(j)}!} \right) \right] + \quad (14) \\ & \sum_{i=1}^k \left[\log \binom{n_i + J_P - 1}{J_P - 1} + \log \left(\frac{n_i!}{n_{i,1}^{(p)}! n_{i,2}^{(p)}! \dots n_{i,J_p}^{(p)}!} \right) \right]. \end{aligned}$$

All parameters and variables in the objective function are explained in previous subsections.

4.3 Algorithm

Once the objective function is established, the next problem is how to find a discretization model that minimizes the objective function. Note that the objective function is cumulative on intervals. Therefore, if a partition of continuous variable X into I intervals with lengths n_1, n_2, \dots, n_k is an optimal discretization policy, then $\{n_2, n_3, \dots, n_k\}$ is optimal for the subproblem, i.e., the last $N - n_1$ instances of the dataset which has been sorted in ascending order of D_X . With this property, we can adapt dynamical programming to solve this optimization problem.

If there is no repeated value in D_X , all middle points of two consecutive values can be discretization edges. If there exists repeated values in D_X , only the middle points of two consecutive and different values can be discretization edges. Assume there are M unique values in D_X , then $D'_X = \{x'_1, x'_2, \dots, x'_M\}$ are the unique values in ascending order. Since D_X is also sorted, let $b = \{b_0, b_1, b_2, \dots, b_M\}$ be an increasing sequence of integers such that $b_0 = 0$ and $x'_i = x_{b_{i-1}+1} = x_{b_{i-1}+2} = \dots = x_{b_i}$. By this definition, the allowable discretization positions are $d_i = (x_{b_{i+1}} + x_{b_i})/2$ for all $i = 1, 2, \dots, M - 1$.

In order to save runtime, we first calculate the following function $g(u, v)$ for a interval I_q starting from x_u to x_v with all u, v satisfying $u \leq v$:

$$\begin{aligned}
g_p(u, v) &= \log \binom{n_q + J_P - 1}{J_P - 1} + \log \left(\frac{n_q!}{n_{q,1}^{(p)}! n_{q,2}^{(p)}! \cdots n_{q,J_p}^{(p)}!} \right) \\
g_c(u, v) &= \sum_{j=1}^{n_c} \sum_{l=1}^{L_j} \left[\log \binom{n_{q,l} + J_j - 1}{J_j - 1} + \log \left(\frac{n_{q,l}!}{n_{q,1,l}^{(j)}! n_{q,2,l}^{(j)}! \cdots n_{q,J_j,l}^{(j)}!} \right) \right] \\
g(u, v) &= g_p(u, v) + g_c(u, v).
\end{aligned} \tag{15}$$

The evaluation of function $g_p(u, v)$ and $g_c(u, v)$ can be done as Algorithm 1 and Algorithm 2, respectively. Algorithm 2 is an illustrative version of calculation $g_c(u, v)$. The more time-efficient but structurally complicated algorithm for calculating $g_c(u, v)$ is shown in Appendix. Notice that due to repeated values in D_X , for some pairs of u and v , $g(u, v)$ might depend on the sorting method. However, this doesn't influence the optimization result, since these pairs of u and v will not form valid intervals. Now we are able to solve the optimization problem with the objective function in Equation 14. The dynamical programming procedure is shown in Algorithm 3. As discussed before, the runtime of Algorithm 3 is $O(n_c v_{max} \cdot N^2 + v_{max}^{n_p} \cdot N^2)$, where v_{max} is the largest cardinality number of variables that directly connects X . There are other methods that lead to suboptimal results but runs faster than dynamical programming method. Please refer to (Boullé 2006).

Algorithm 1 Calculation of function $g_p(u, v)$ for all $u \leq v$

```

1: Initialize  $g_p$  as an  $N$  by  $N$  matrix that all elements are 0.
2:  $count_p$  is an  $N$  by  $N$  by  $||I_X||$  matrix such that  $count_p(u, v, w)$  is the number of instances
   from  $x_u$  to  $x_v$  with value of  $I_X$ . This matrix can be calculated in  $O(||I_X|| \cdot N^2)$ 
3: for  $u = 1$  to  $N$  do
4:   for  $v = u$  to  $N$  do
5:      $g_p(u, v) \leftarrow g_p(u, v) + \log((v - u + J_p)!) - \log((J_p - 1)!)$ 
6:     for  $w = 1$  to  $||I_X||$  do
7:        $g_p(u, v) \leftarrow g_p(u, v) - \log(count_p(u, v, w)!)$ 
8:     end for
9:   end for
10: end for

```

4.4 Approximation

Notice that Algorithm1 leads to the non-polynomial runtime $O(v_{max}^{n_p} \cdot N^2)$ in the discretization process. Therefore, we have an empirical approximation to the objective function that can reduce runtime and still preserve the quality of discretization. As shown in Section 6, the approximated algorithm is more

Algorithm 2 Calculation of function $g_c(u, v)$ for all $u \leq v$

```

1: for  $j = 1$  to  $n_c$  do
2:    $count_j$  is an N by N by  $\|S_j\|$  by  $\|C_j\|$  matrix such that  $count_j(u, v, q, r)$  is the
   number of instances from  $x_u$  to  $x_v$  with  $q$ th value of  $S_j$  and  $r$ th value of  $C_j$ .
3:   for  $u = 1$  to  $N$  do
4:     for  $v = u$  to  $N$  do
5:       for  $q = 1$  to  $\|S_j\|$  do
6:          $n' \leftarrow \sum_{r=1}^{J_j} count_j(u, v, q, r)$ 
7:          $g(u, v) \leftarrow g(u, v) + \log((v - u + 1 + n')!) - \log((v - u + 1)!)
8:         for  $r = 1$  to  $J_j$  do
9:            $g(u, v) \leftarrow g(u, v) - \log((count_j(u, v, q, r))!)$ 
10:        end for
11:      end for
12:    end for
13:  end for
14: end for$ 
```

Algorithm 3 Discretization of one continuous variable

```

1: Let  $S(u)$  be the optimal objective function value of a subproblem starting from instance
   1 to instance  $u$ .
2: Let  $Disc(u)$  be the optimal discretization of a subproblem starting from instance 1 to
   instance  $u$ .
3: Let  $L(u) = -\log(1 - \exp(-M \cdot \frac{x_{b(u)+1} - x_{b(u)}}{x_N - x_1}))$  for  $u = 1, 2, \dots, M - 1$  and  $L(M) = 0$ .
4: Calculate  $g(u, v)$  for all  $u \leq v$  as Algorithm 1.
5: for  $v = 1$  to  $M$  do
6:   if  $v = 1$  then
7:      $S(v) = g(1, b(v)) + L(v)$ 
8:      $Disc(v) = \{(x_{b(v)} + x_{b(v)+1})/2\}$ 
9:   else
10:     $s = \infty$  and  $boundary = \infty$ 
11:    for  $u = 1$  to  $v$  do
12:       $s' \leftarrow S(u) + g(b(u) + 1, b(v)) + \frac{M(x_{b(v)} - x_{b(u)+1})}{x_N - x_1} + L(v)$ 
13:      if  $s' < s$  then
14:         $s \leftarrow s'$ 
15:         $boundary = (x_{b(u)} + x_{b(u)+1})/2$ 
16:      end if
17:    end for
18:     $S(v) = s$ 
19:     $Disc(v) = Disc(u) \cup \{boundary\}$ 
20:  end if
21: end for

```

sensitive to the distribution of other variables' values and usually leads to slightly more discretization edges. With the approximation, we can reduce runtime of Algorithm 1 from $O(v_{max}^{n_p} \cdot N^2)$ to $O(v_{max} \cdot n_p \cdot N^2)$. The rigorous proof of the approximation have not been proposed.

For the dominator of the last factor in Equation 11, we have the following approximation

$$\frac{n_i!}{n^{(p)}_{i,1}! n^{(p)}_{i,2}! \dots n^{(p)}_{i,J_P}!} \approx \prod_{q=1}^{n_p} \frac{n_i!}{n^{(p_q)}_{i,1}! n^{(p_q)}_{i,2}! \dots n^{(p_q)}_{i,J_{p_q}}!}, \quad (16)$$

where $J_{p_q} = ||P_q||$ and $n_{i,j}^{(p_q)}!$ is the number of instances in i th interval with j th value of P_q . Then the calculation of corresponding approximated $g_p(u, v)$ can be done as Algorithm 4. With the approximation, the approximate objective function is

$$\begin{aligned} & \sum_{i=1}^{k-1} -\log(1 - \exp(-M \cdot \frac{x_{s_i+1} - x_{s_i}}{x_N - x_1})) + \sum_{i=1}^k M \cdot \frac{x_{s_i} - x_{s_{i-1}+1}}{x_N - x_1} + \\ & \sum_{j=1}^{n_c} \sum_{i=1}^k \sum_{l=1}^{L_j} \left[\log \binom{n_{i,l} + J_j - 1}{J_j - 1} + \log \left(\frac{n_{i,l}!}{n_{i,1,l}^{(j)}! n_{i,2,l}^{(j)}! \dots n_{i,J_j,l}^{(j)}!} \right) \right] + \quad (17) \\ & \sum_{i=1}^k \left[\log \binom{n_i + J_P - 1}{J_P - 1} + \sum_{q=1}^{n_p} \log \left(\frac{n_i!}{n_{i,1}^{(p_q)}! n_{i,2}^{(p_q)}! \dots n_{i,J_{p_q}}^{(p_q)}!} \right) \right]. \end{aligned}$$

Algorithm 4 Calculation of approximated function $g_p(u, v)$ for all $u \leq v$

```

1: Initialize  $g_p$  as an  $N$  by  $N$  matrix that all elements are 0.
2: for  $q = 1$  to  $n_p$  do
3:    $count_{p_q}$  is an  $N$  by  $N$  by  $||P_q||$  matrix such that  $count_{p_q}(u, v, w)$  is the number
   of instances from  $x_u$  to  $x_v$  with  $w$ th value of  $P_q$ . This matrix can be calculated in
    $O(||P_q|| \cdot N^2)$ 
4:   for  $u = 1$  to  $N$  do
5:     for  $v = u$  to  $N$  do
6:        $g_p(u, v) \leftarrow g_p(u, v) + \log((v - u + ||P_q||)!) - \log((||P_q|| - 1)!) +$ 
7:       for  $w = 1$  to  $||P_q||$  do
8:          $g_p(u, v) \leftarrow g_p(u, v) - \log(count_{p_q}(u, v, w)!) +$ 
9:          $\log((v - u + ||P_q|| - w + 1)!) - \log((||P_q|| - w)!) +$ 
10:         $\log(w!)$ 
11:       end for
12:     end for
13:   end for

```

With Equation 17, we have a more clear way to see that how child variables and parent variables contribute to the objective function differently. For example, in the left graph of Figure 2, the corresponding square bracket in Equation 17 is

$$\begin{aligned} & \sum_{i=1}^k \left\{ \left[\log \binom{n_i + J_{C_1} - 1}{J_{C_1} - 1} + \log \left(\frac{n_i!}{n_{i,1,\emptyset}^{(1)}! \dots n_{i,J_{C_1},\emptyset}^{(1)}!} \right) + \right. \right. \\ & \quad \left. \log \binom{n_i + J_{C_2} - 1}{J_{C_2} - 1} + \log \left(\frac{n_i!}{n_{i,1,\emptyset}^{(2)}! \dots n_{i,J_{C_2},\emptyset}^{(2)}!} \right) \right] + \quad (18) \\ & \quad \left. \left[\log \binom{n_i + J_P - 1}{J_P - 1} + \log \left(\frac{n_i!}{n_{i,1}^{(p_1)}! \dots n_{i,J_{p_1}}^{(p_1)}!} \right) + \frac{n_i!}{n_{i,1}^{(p_2)}! \dots n_{i,J_{p_2}}^{(p_2)}!} \right] \right\} \end{aligned}$$

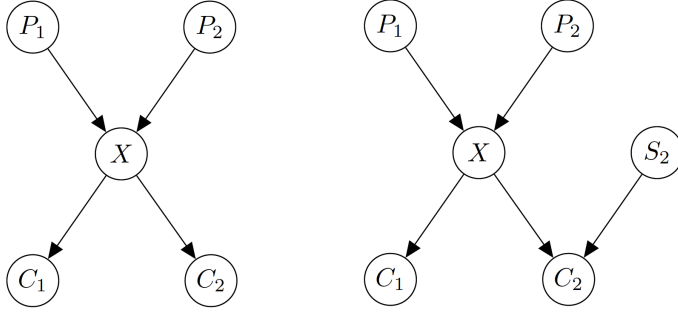


Fig. 2 Two example networks

In Equation 18, each child variable carries two terms as the first square bracket, and two parent variables only carry three terms as the second square bracket. Therefore, in this case, child variables are more determinative than parent variables, even though the number of child variables is same as the number of parent variables. However, if C_2 has another parent variables, as shown in the right graph of Figure 2, the importances of C_2 will be debilitated, since the information from C_2 is now adulterated by S_2 . These arguments show that our proposed method, either before or after approximation, indeed involves graph structures to discretize continuous variables.

5 Discretizing Multiple Continuous Variable with Structure Learning

5.1 Discretization of Multiple Continuous Variables

If there are multiple continuous variables in a Bayesian network, we iterate the one-variable discretization method discussed above for all continuous variables. While discretizing a continuous variable, other continuous variables must be discretized, either by prediscretization before the iteration starts or by the discretization results of previous iteration. The prediscretization here can be done by equal-width discretization, which is defined as follows:

$$M_X = \{min_X, min_X + \delta, min_X + 2\delta, \dots, \delta, max_X\}, \quad (19)$$

where min_X and max_X are minimal and maximal values of D_X , respectively, $\delta = (max_X - min_X)/k$ and k is the desired number of intervals after equal-width discretization. k is same for all continuous variables and is set to be median value of all discrete variables' cardinalities. After prediscretization, we iterate the one-variable discretization on each continuous variable in the following order: from the continuous variable with highest topological order (leaves) to the continuous variable with lowest topological order (root). We call a series of the iterations that follows this order one time as a cycle. The

advantage of the order is that, for the first cycle of iteration, we can use less number of unsupervised discretization result. For example, in the right graph of Figure 2, assume S_2 is the only discrete variable. If we begin the iteration with P_1 , then the discretization of P_1 will involve the prediscrretization results of P_2 and X . However, if we begin the iteration with C_1 , then the discretization of C_1 will only involve the prediscrretization result of X .

We stop the iterations and output the discretization result on each continuous variable when the number of intervals and the positions of discretization edges converge. Since the convergence is not guaranteed, we also set up a maximal number of cycles to prevent infinite iterations. In our test on real-world data, the iteration results usually converge within few cycles. Even if it does not converge, 10 cycles usually produce good enough discretization result. The summary of multi-variable discretization is shown as Algorithm 5.

Algorithm 5 Discretization of multiple continuous variables

```

1: Let  $Disc(i)$  is the discretization on  $i$ th continuous variable. If variable  $i$  is discrete,
    $Disc(i)$  contains no information.
2: Let  $D_i^*$  is the discrete data of  $i$ th variable. If variable  $i$  is continuous, then  $D_i^*$  is
   discretized  $D_i$  by  $Disc(i)$ .
3: Let  $C$  be the set of continuous variables in a reverse topological order.
4: Let  $n$  be the number of variables in the given Bayesian network,  $k$  is the largest cardinality
   number of discrete variables.
5: % Prediscrretize continuous variables,  $bound_{cycle}$  is the upper bound of cycles.
6: for  $i = 1$  to  $n$  do
7:   if  $i \in C$  then
8:      $Disc(i) =$  prediscrretization result of  $k$  intervals
9:   end if
10: end for
11: % Iterate one-variable discretization until converge
12:  $cycle = 0$ 
13: while  $Disc$  is not converged &  $cycle \leq bound_{cycle}$  do
14:    $cycle \leftarrow cycle + 1$ 
15:   for  $j = 1$  to  $\|C\|$  do
16:     Let  $D_{\setminus C(j)}^*$  be the discrete data of all variables except  $D_{C(j)}^*$ .
17:      $Disc(C(j)) =$  one-variable discretization result with  $D_{C(j)}$  and  $D_{\setminus C(j)}^*$ .
18:      $D_{C(j)}^* =$  the discretized data of  $D_{C(j)}$  by latest discretization policy  $Disc(C(j))$ 
19:   end for
20: end while

```

5.2 Discretization of Continuous Variables While Structure Learning

6 Experiments

References

M. Boullé. Modl: A bayes optimal discretization method for continuous attributes. *Machine Learning*, pages pp131–165, 2006.

- G. F. Cooper and E. Herskovits. A bayes method for the induction of probabilistic network from data. *Machine Learning*, 9:pp.309–147, 1992.
- T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stain. *Introduction to algorithms*. The MIT Press, 3 edition, 2009.
- J. Dougherty, R. Kohavi, and M. Sahami. Supervised and unsupervised discretization of continuous features. In *12th International Conference on Maching Learning (ICML)*. Morgan Kaufman Publishers, San Francisco, CA, 1995.
- U.M. Fayyad and K.B. Irani. Multi-interval discretization of continuous-valued attributes for classification learning. In *13th International Joint Conference on Artificial Intelligence*, pages pp.1022–1027. Morgan Kaufman, 1993.
- N. Friedman and M. Goldszmidt. Discretizing continuous attributes while learning bayesian networks. In *13th International Conference on Maching Learning (ICML)*. Morgan Kaufman Publishers, San Francisco, CA, 1996.
- N. Friedman, D. Geiger, and M. Goldszmidt. Bayesian network classifiers. *Machine Learning*, 29:pp.131–163, 1997.
- P. D. Grünwald. *Minimum description length principle*. The MIT Press, 2007.
- R.C. Holte. Very simple classification rules perform well on most commonly used datasets. *Machine Learning*, 11:pp.63–90, 1993.
- A. Ihler and D. McAllester. Particle belief propagation. In *Artificial Intelligence and Statistics*, 2009.
- R. Kerber. Chimerge: Discretization of numeric attributes. In *10th National Conference on Artificial Intelligence*, pages pp.123–128. MIT Press, 1992.
- M. J. Kochenderfer. *Decision making under uncertainty*. MIT Lincoln Laboratory Series, 2015.
- D. Koller and N. Friedman. *Probabilistic graphical models: principles and techniques*. The MIT Press, 2009.
- A.V. Kozlov and D. Koller. Nonuniform dynamic discretization in hybrid networks. In *13th International Conference on Uncertainty in Artificial Intelligence (UAI)*, 1997.
- J. L. Lustgarten, S. Viswewaran, Gopalakrishnan V., and G. F. Cooper. Application of an efficient bayesian discretization method to biomedical data. *BMC Bioinformatics*, 2011.
- S. Monti and G.F. Cooper. A multivariate discretization method for learning bayesian networks from mixed data. In *14th International Conference on Uncertainty in Artificial Intelligence (UAI)*, 1998.
- J. Pearl. Probabilistic reasoning in intelligent systems: networks of plausible inference. *Morgan Kaufman Publishers, Inc.*, 1988.
- J. Rissanen. Modeling by shortest data description. *Automatica*, pages pp.465–471, 1978.
- H. Stech and T. Jaakkola. Predictive discretization during model selection. In *11th International Conference on Artificial Intelligence and Statistics*, 2007.
- P. M. B. Vitényi and M. Li. Minimum description length induction, bayesianism, and kolmogorov complexity. *IEEE Transactions on Information Theory*, 2000.

Y. Weiss and W. T. Freeman. Correctness of belief propagation in gaussian graphical models of arbitrary topology. *Neural Computation*, 2001.

Appendices