

Efficient and Accurate Learning of Bayesian Networks using Chi-Squared Independence Tests

Yi Tang and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition
 yitang@buffalo.edu, srihari@cedar.buffalo.edu

Abstract

Bayesian network structure learning is a well-known NP-complete problem, whose solution is of importance in machine learning. Two algorithms are proposed, both of which assess dependency between variables using the chi-squared test of independence between pairs of variables and the log-likelihood evaluation criterion for the network. The first determines the effect of adding a potential edge (in both directions) on the log-likelihood. The second uses K-L divergence to determine direction, and edges to be included are determined by thresholding normalized chi-squared statistics. Experiments on multinomial data show that the proposed algorithms are more efficient and accurate than an optimized branch and bound algorithm, and human experts.

1. Introduction

Learning the structure of Bayesian networks (BNs) [1, 6, 7, 10] that approximate the joint distribution of a large number of variables is an important problem in machine learning and data mining. Existing methods for structure learning [7] in BNs can be divided into three types: *constraint* based, *score* based and *Bayesian model averaging* methods. By viewing the BN as a representation of dependencies, constraint based methods attempt to find a network structure that best explains dependencies. But it is sensitive to errors in testing single dependencies [7]. Score based methods view learning as a *model selection* problem; by defining a scoring function which assesses the fitness of each model, it searches for a high-scoring network structure. As the search space is super-exponential, enumerating scores for all models is often *NP-hard*. Therefore it has to resort to heuristic search. Examples are the K2 algorithm [2], and the optimized branch and bound algorithm [4]. The third type of methods, Bayesian model averaging,

make predictions by averaging across all possible structures. But the disadvantage is that some may not have closed form. Recently Peters *et al.* [8] proposed an algorithm that infers the causal structure of discrete variables using additive noise models. It also points out that the limitation of the χ^2 test on small data sets.

We propose two new algorithms which score each candidate network structure based on the standard likelihood function. Instead of exhaustively searching the entire solution space, they utilize the constraints obtained from statistic χ^2 tests [12] as guidance to *intelligently* and *incrementally* search for the structure with a higher score until reaching the *optimal*. Two algorithms differ in how they evaluate the goodness of fit of each candidate structure, i.e. the first algorithm uses the global criteria to determine directionality whereas the second algorithm decides the direction using local criteria. The two algorithms overcome the limitation of χ^2 tests by using a majority voting procedure (in Algorithm 1) and defining a normalized χ^2 statistic measure (in Algorithm 2), i.e. the optimality of one network over the others is *not* based on a single test but multiple tests.

2. Learning Bayesian Networks

Basically, a BN is a collection of *directed dependencies* (or *causality*) between variables. The goal of our algorithm is to find a structure that gives the best representation of all directed dependencies between variables. The goodness of a BN is based on its capability to generalize to unseen samples in terms of the accuracy in specific tasks.

We determine the degree of dependency between each pair of random variables using Pearson's χ^2 test [12]. The original χ^2 test computes the χ^2 statistic, then accepts or rejects the null hypothesis by determining whether the estimated *p*-value exceeds a pre-defined significance level or not. However, this method does not work well in practical applications as data contain

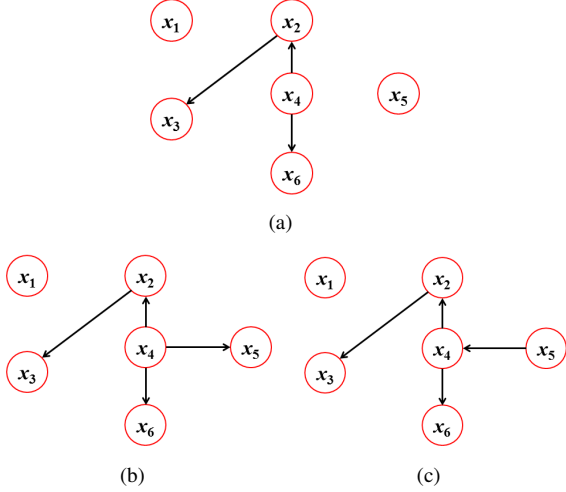


Figure 1: Global algorithm. Edge (x_4, x_5) is being considered. (a) the best Bayesian network (BN) G^* before considering edge (x_4, x_5) ; (b) candidate network G_{c_1} with edge $(x_4 \rightarrow x_5)$; (c) candidate network G_{c_2} with edge $(x_5 \rightarrow x_4)$.

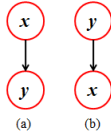


Figure 2: Local algorithm.

noise. Instead, we need to either find an additional criterion base on which the algorithm can make a decision, or learn an optimal threshold for χ^2 statistic, if it exists. Based on whether an explicit threshold for χ^2 statistic has to be learned, we have designed two algorithms for BN structure learning. Algorithm 1 examines one edge (x_i, x_j) at a time, and determines whether or not the edge should be added and which direction should be used by comparing the three Bayesian networks (BNs): the best BN before considering the edge (G^*), the BN with $(x_i \rightarrow x_j)$ added (G_{c_1}), and the BN with $(x_j \rightarrow x_i)$ added (G_{c_2}), as shown in Fig. 1 (a), (b) and (c) respectively. We evaluate the goodness of each BN by the (log) likelihood function of the data, and choose the one with the highest log likelihood score as updated G^* , and then proceed to the next edge.

Similar to Algorithm 1, Algorithm 2 also performs the Chi-squared test between all pairs of variables. To determine the direction of each edge (x, y) , it considers the joint distribution $P(x, y)$ and its two possible factorizations, namely $Q_1(x, y) = P(x)P(y|x)$ and $Q_2(x, y) = P(y)P(x|y)$, and computes the K-L divergences between them. The factorization with a lower value of K-L divergence indicates the direction, i.e. if $\text{KL}(P||Q_1) < \text{KL}(P||Q_2)$, then the direction should be $(x \rightarrow y)$, else it should be $(y \rightarrow x)$, as shown in Fig. 2.

2.1. Preliminaries

Independence Test. A suitable test for independence of two categorical variables is the *Pearson's Chi-Squared Test*. It is a goodness-of-fit test in which the data are compared with a model based on the null hypothesis of independence[5]. Let X and Y be two multinomial variables governed by distributions $P(X = x_i)$ and $P(Y = y_j)$. Given a data set \mathcal{D} with total m samples, $O[x_i, y_j]$ is *observed* count for each *joint assignment* of $X = x_i$ and $Y = y_j$. Given the null hypothesis that X and Y are independent, *expected* count for (x_i, y_j) is $E[x_i, y_j] = mP(x_i)P(y_j)$. Then the χ^2 statistic [12] for data set \mathcal{D} is defined by

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

Parameter Estimation. For a multinomial variable x with k possible values $\{x_1, \dots, x_k\}$, we place a uniform Dirichlet prior on its parameters $\theta = \{\theta_1, \dots, \theta_k\}$, i.e. $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$, where $\alpha_1 = \dots = \alpha_k = 1$. Now consider a data set \mathcal{D} of independent observations where the number of observations of $x = x_i$ ($i = 1, \dots, k$) is o_i . Then the *sufficient statistics* $O = \{o_1, \dots, o_k\}$ is distributed by $\{o_1, \dots, o_k\} \sim \text{Multinomial}(\theta_1, \dots, \theta_k)$. Because of the conjugacy between Dirichlet distribution and multinomial distribution, the posterior distribution of θ given O , is also Dirichlet, $\theta|O \sim \text{Dirichlet}(\alpha'_1, \dots, \alpha'_k)$, where $\alpha'_i = \alpha_i + o_i, \forall i = 1, \dots, k$.

Score. The log-likelihood of a data set \mathcal{D} with N samples (i.i.d.) for a BN G with n variables, is given by

$$s(\mathcal{D}|G) = \sum_{i=1}^N \sum_{j=1}^n \log P(x_{ij}|\text{pa}_{ij}), \quad (2)$$

where x_{ij} is the j^{th} variable (and its value) in the i^{th} sample, pa_{ij} is the parent(s) of the variable x_{ij} .

2.2. Global Algorithm

This algorithm first computes χ^2 statistics between all pairs of nodes, and constructs a set E_p of node pairs in *non-increasing* order of χ^2 statistics. Starting from the simplest network $G_c = (V, E, \theta)$ in which $E = \emptyset$, the algorithm selects the node pair $(v_i, v_j) \in E_p$ with maximum χ^2 statistic, sets up two directed edges $(v_i \rightarrow v_j)$ and $(v_j \rightarrow v_i)$. By adding each of them into the original graph G_c , it creates two candidate graphs $G_{c_1} = (V, E_{c_1}, \theta_{c_1})$ and $G_{c_2} = (V, E_{c_2}, \theta_{c_2})$, in which $E_{c_1} = E + \{v_i \rightarrow v_j\}$, and $E_{c_2} = E + \{v_j \rightarrow v_i\}$. The parameters of G_c, G_{c_1} , and G_{c_2} are estimated as in Section 2.1. To evaluate the suitability of these graphs, cross validation is used. Basically for each run, the

two thirds of samples are randomly selected from the data and used for training, and the rest is used for testing. The goodness of each network is measured by its log-likelihood score s that is defined in Eq. 2. The algorithm will choose the best model G^* so far with the highest score s among (G_c, G_{c_1}, G_{c_2}) , set current graph $G_c = G^*$, remove (v_i, v_j) out of E_p . Then it repeats the previous procedures, until all node pairs have been examined, i.e. $E_p = \emptyset$.

Algorithm 1: Learning w/o thresholding

Input: set $S = \{v_1, \dots, v_n\}$ of n random variables, training data \mathcal{D}_r and validation data \mathcal{D}_v

Output: an optimal graph $G^* = \{V, E^*, \theta^*\}$

1. Compute pairwise χ^2 statistics, $\chi^2(v_i, v_j)$
2. Construct ordered set $E_p = \{(v_i, v_j) | i \neq j\}$ in descending order of χ^2
3. Set $G^* = \{V, E^*, \theta^*\}$, where $V = S, E^* = \emptyset$
4. Learn θ^* using \mathcal{D}_r
5. Set $k = 1$
6. **Repeat**
7. Pick the first pair (v_i, v_j) from E_p
8. Create $G_{c_1} = (V, E_{c_1}, \theta_{c_1}), E_{c_1} = E^* + \{v_i \rightarrow v_j\}$
9. **If** isDAG(G_{c_1})
10. **Then** Estimate θ_{c_1} from \mathcal{D}_r
11. Compute s_{k-1}, s_{c_1} for G^* and G_{c_1}
12. **If** $s_{c_1} > s_{k-1}$
13. **Then** $G^* = G_{c_1}$
14. Create $G_{c_2} = (V, E_{c_2}, \theta_{c_2}), E_{c_2} = E_{c_1} - \{v_i \rightarrow v_j\} + \{v_j \rightarrow v_i\}$
15. **If** isDAG(G_{c_2})
16. **Then** Estimate θ_{c_2} from \mathcal{D}_r
17. Compute s_{k-1}, s_{c_2} for G^* and G_{c_2}
18. **If** $s_{c_2} > s_{k-1}$
19. **Then** $G^* = G_{c_2}$
20. Increment k by 1
21. $E_p = E_p - \{(v_i, v_j)\}$
22. **Until** $E_p = \emptyset$
23. **return** G^*

The function isDAG(G) returns a Boolean variable that indicates whether the Bayesian graph G is a *directed acyclic graph* (DAG) or not. In Algorithm 1, step 1 takes time $O(n^2)$, step 2 takes $O(n^2 \log(n))$, step 4 takes time $O(n \cdot 2^d)$, where d is the maximum number of parents per node, and the time required in step 6 through 22 is bounded by $O(n^2)$, so the running time of Algorithm 1 is $O(n^2 \log(n))$.

2.3. Local Algorithm

The *Kullback-Leibler divergence* is used to determine dependency direction. Given two variables (x, y) and two possible decompositions of $P(x, y)$, i.e. $Q_1 = P(x)P(y|x)$ and $Q_2 = P(y)P(x|y)$, we can compute

$KL(P||Q_1)$ and $KL(P||Q_2)$, and choose the decomposition with lower divergence. To control graph structure complexity, we introduce a threshold t for degree of dependency, so that only those edges with χ^2 statistic above t will be taken into account. To find the optimal value of t , we divide the data into 3 subsets for training, validation, and testing. Suppose that there is a range $[t_l \ t_u]$ for valid thresholds. For each possible threshold $t \in [t_l \ t_u]$, we construct a BN using those directed edges (avoiding forming loops) with the χ^2 statistic t to learn the parameters with training data, and then assess its performance with validation data. The value of t which leads to the best performance on validation data set is regarded as the optimal. We use a priority queue Q to store all the possible edges, in which each entry is a tuple (e_k, d_k) ($k = 1, \dots, \frac{n(n-1)}{2}$) where e_k is an undirected edge (v_i, v_j) or directed edge $(v_i \rightarrow v_j)$, d_k is the degree of dependency of the edge e_k computed by Eq. 3. All the entries are sorted in non-decreasing order of d_k . Below we first define normalized χ^2 statistics, then summarize our algorithm in Algorithm 2.

Normalized χ^2 statistics. The χ^2 statistic defined in Eq. 1 is *sensitive* to the number of samples m because its value increases as m increases, even if the distribution stays the same. We define the normalized χ^2 statistic by

$$\chi_n^2(\mathcal{D}) = \sum_{i,j} \frac{(O[x_i, y_j]/m - E[x_i, y_j]/m)^2}{E[x_i, y_j]/m}, \quad (3)$$

which truly represents the distribution of the data, irrespective of the number of observations for each joint assignment.

Algorithm 2: Learning with thresholding

Input: set $S = \{v_1, \dots, v_n\}$ of n random variables (r.v.), training data \mathcal{D}_r and validation data \mathcal{D}_v

Output: an optimal graph $G^* = \{V, E^*, \theta^*\}$

1. Compute pairwise normalized χ^2 statistics
2. Construct a priority queue $Q = \{(e_k, d_k)\}$
3. **For** each entry $(e_k, d_k) \in Q$ **Do**
4. Determine its directionality using MI
5. **For** $t := t_u$ **downto** t_l **Do**
6. Create $G_c = \{S, E_c, \theta_c\}, E_c = \{e_k | d_k \geq t\}$
7. Learn θ_c from \mathcal{D}_r
8. Compute score s_c for G_c using \mathcal{D}_v
9. $t = t - \eta$ /* η is learning rate.*/
10. Find the optimal value t^* by plotting s_c against t
11. Create $G^* = \{V, E^*, \theta^*\}, E^* = \{e_k | d_k \geq t^*\}$
12. Learn θ^* from \mathcal{D}_r
13. **return** G^*

3. Evaluation

Two data sets were used: (i) handwritten "th" described by six features[9], with 3,125 samples of

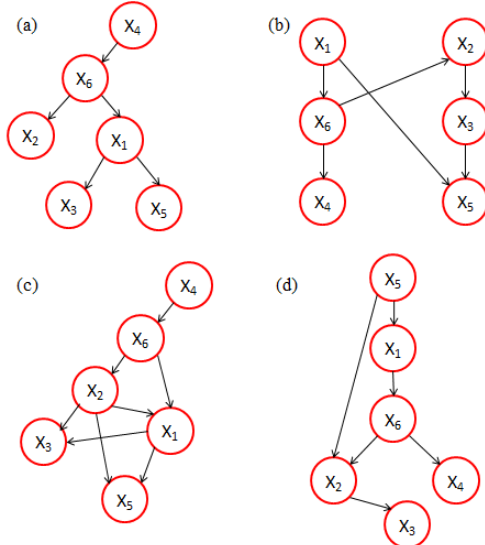


Figure 3: Bayesian networks learned from “th” data. (a) B & B Algorithm; (b) Human designed based on domain knowledge; (c) Algorithm 1; (d) Algorithm 2.

“th” from 1, 254 documents written by 499 writers, and (ii) “and” described by nine features[11], with 10, 111 samples written by 1, 555 authors, in which two writing styles including cursive and hand-print are differentiated. There is one feature vector per document, i.e., all the “th”s in one document represented by a single feature vector. A branch and bound algorithm (B & B) [4, 3] claiming “global optimality” and human-designed BNs were used for baseline comparison.

Table 1: Comparison of Algorithms & Human for “th”.

Methods	Log-likelihood s	Avg. Learn Time (sec.)	Sensitivity σ_s
B & B	-2598.4	1.8	25.7
Human Design	-2595.8	N/A	-
Alg. 1	-2588.9	1.6	20.0
Alg. 2	-2605.3	1.0	21.4

Table 2: Comparison for cursive “and”.

Methods	Log-likelihood s	Avg. Learn Time (sec.)	Sensitivity σ_s
B & B	-4008.5	2.1	36.0
Human Design	-4136.7	N/A	39.5
Alg. 1	-3957.5	1.8	35.9
Alg. 2	-3983.4	1.4	36.0

To study algorithm sensitivity, we ran validation tests 1000 times for all candidate graphs on *randomized* data. For each run, we randomly selected $\frac{2}{3}$ of the whole data for training, and the rest for test. Three evaluation metrics were: log-likelihood of the test data, s , average learning time \bar{T} , and sensitivity to training data, σ_s , defined by the *standard deviation* of the log-likelihood scores from multiple runs. BNs learned from “th” data are shown in Fig. 3. Their performances are in Tables 1-3. Both algorithms are more accurate and efficient than B & B and human design. Between the two, the first is more accurate (as it uses global criteria) and the second is more efficient (because of thresholding).

Table 3: Comparison for *handprinted* “and”.

Methods	Log-likelihood s	Avg. Learn Time (sec.)	Sensitivity σ_s
B & B	-1416.9	1.6	34.4
Human Design	-1494.8	N/A	40.8
Alg. 1	-1401.4	1.3	33.0
Alg. 2	-1405.7	1.1	33.7

4. Conclusion

Two algorithms for BN structure learning have been proposed, with the goal of maximizing the likelihood of the data. Exponential search is reduced to polynomial-time by exploiting dependency constraints between variables. Algorithm 1 examines *each* edge without thresholding. Algorithm 2 reduces the number of potential edges by determining the direction of each edge using K-L divergence, and learns the optimal threshold for degree of dependency. The two algorithms performed better than a state-of-the-art branch-and-bound algorithm and human experts.

References

- [1] C. K. Chow and C. N. Liu. Approximating discrete probability distributions with dependence trees. *IEEE Trans. on Information Theory*, 14(3):462–467, 1968.
- [2] G. Cooper and E. Herskovits. A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, 9(4):309–347, 1992.
- [3] C. P. de Campos and Q. Ji. Efficient structure learning of Bayesian networks using constraints. *Journal of Machine Learning Research*, 12:663–689, 2011.
- [4] C. P. de Campos, Z. Zeng, and Q. Ji. Structure learning of Bayesian networks using constraints. In *Proc. of ICML 2009*, pages 113–120, 2009.
- [5] D. Hand, H. Mannilla, and P. Smyth. *Principles of Data Mining*. MIT Press, 2001.
- [6] S. Kirshner, P. Smyth, and A. W. Robertson. Conditional Chow-Liu tree structures for modeling discrete-valued vector time series. In *Proc. of the 20th conference on UAI*, pages 317–324, 2004.
- [7] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [8] J. Peters, D. Janzing, and B. Schölkopf. Causal inference on discrete data using additive noise models. *IEEE TPAMI*, 33(12):2436–2450, Dec. 2011.
- [9] S. N. Srihari. Evaluating the rarity of handwriting formations. In *Proc. of ICDAR*, pages 618–622, 2011.
- [10] J. Suzuki. A generalization of the Chow-Liu algorithm and its application to statistical learning. In *Proc. of Intl. Conf. on Artificial Intelligence*, pages 478–486, 2010.
- [11] Y. Tang and S. N. Srihari. Learning Bayesian networks for likelihood ratio computation. In *Proc. IWCF*, 2012.
- [12] L. Wasserman. *All of Statistics: A Concise Course in Statistical Inference*. Springer, 2003.