

# GlobalMIT: Learning Globally Optimal Dynamic Bayesian Network with the Mutual Information Test (MIT) Criterion

Nguyen Xuan Vinh<sup>1</sup>, Madhu Chetty<sup>1</sup>, Ross Coppel<sup>2</sup> and Pramod P. Wangikar<sup>3</sup>

<sup>1</sup>Gippsland School of Information Technology, Monash University, Victoria, Australia,

<sup>2</sup>Faculty of Medicine, Nursing and Health Sciences, Monash University, Victoria, Australia,

<sup>3</sup>Chemical Engineering Department, Indian Institute of Technology Bombay, India.

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: Trey Ideker

## ABSTRACT

**Motivation:** Dynamic Bayesian networks (DBN) are widely applied in modeling various biological networks including the gene regulatory network (GRN). Due to the NP-hard nature of learning static Bayesian network structure, most methods for learning DBN also employ either local search such as hill-climbing, or a meta stochastic global optimization framework such as genetic algorithm or simulated annealing.

**Results:** This paper presents GlobalMIT, a toolbox for learning the globally optimal DBN structure from gene expression data. We propose using a recently introduced information theoretic based scoring metric named mutual information test (MIT). With MIT, the task of learning the globally optimal DBN is efficiently achieved in polynomial time.

**Availability:** The toolbox, implemented in Matlab and C++, is available at <http://code.google.com/p/globalmit>.

**Contact:** {vinh.nguyen,madhu.chetty}@monash.edu

## 1 INTRODUCTION

Bayesian network (BN) has found applications in modeling various biological networks including the gene regulatory network (GRN). The two important limitations when applying static BN to these domain problems are: (i) BN does not have a mechanism for exploiting the temporal aspect of time-series data, such as time-series microarray data; and (ii) BN does not allow the modeling of cyclic phenomena, such as feed back loops, which are prevalent in biological systems (Yu *et al.*, 2004). These drawbacks have motivated the development of the so-called dynamic Bayesian network (DBN). Its simplest model, the first-order Markov stationary DBN, assumes that both the structure of the network and the parameters characterizing it remain unchanged over time. The value of a variable at time ( $t$ ) is assumed to depend only on the value of its parents at time ( $t - 1$ ). DBN not only accounts for the temporal aspect of time-series data (i.e., an inter time-slice edge must always be directed forward in time), but it also allows the modeling of feedback loops. Since its inception, DBN has received particular interest from the bioinformatics community (Murphy and Mian, 1999; Perrin *et al.*, 2003; Kim *et al.*, 2003; Husmeier, 2003; Yu *et al.*, 2004; Zou and Conzen, 2005; Wilczynski and Dojer, 2009).

## 2 METHOD

Most algorithms to date for learning DBN structure employ a local search strategy, such as hill climbing with random restart, or a meta stochastic global optimization framework such as genetic algorithm or simulated annealing, as exemplified by several softwares such as BANJO (Smith *et al.*, 2006) or bnlearn (Scutari, 2010). This is due to several NP-hardness results in learning static BN structure (see, e.g., Chickering, 1996). Recently, Dojer (2006) has shown otherwise that learning DBN structure, as opposed to static BN, does not necessarily have to be NP-hard. In particular, it was shown that, under some mild assumptions, there are algorithms using the Minimum Description Length (MDL) and BDe scores, which find the globally optimal network with a polynomial worst-case time complexity. These algorithms have been realized within the BNFinder software (Wilczynski and Dojer, 2009). In our experiments, we observed that BNFinder+MDL is very fast, whereas BNFinder+BDe is very time demanding: a single run on a data set of 20 genes and 300 observations can take up to a day (Vinh *et al.*, 2011). This is in concordance with the theoretical worst-case complexity analysis, where the algorithm would have to exhaustively evaluate all possible parent sets of cardinality from 0 to  $p^* - 1$ . Let  $k$  be the number of discrete states of each variable,  $N$  is the number of experiments, then for MDL,  $p_{MDL}^*$  is given by  $\lceil \log_k N \rceil$ , while for BDe,  $p_{BDe}^* = \lceil N \log_{\gamma-1} k \rceil$ , where  $0 < \gamma < 1$  is the network complexity penalty parameter (default value  $\log \gamma^{-1} = 1$  for BNFinder). In general,  $p_{BDe}^*$  scales linearly with the number of data items  $N$ , making its value of less practical interest, even for very small data sets. Although being more expensive, BNFinder+BDe is still recommended over BNFinder+MDL, “due to its exactness in the statistical interpretation” (Wilczynski and Dojer, 2009). Further, de Campos (2006) also showed that BDe seems to learn more accurate networks than MDL (which is also equivalent to the Bayesian Information Criterion—BIC).

Mutual Information Test (MIT) is a recently introduced scoring metric for learning BN (de Campos, 2006). To understand MIT, let  $\mathbf{X} = \{X_1, \dots, X_n\}$  denote the set of  $n$  variables with corresponding  $\{r_1, \dots, r_n\}$  discrete states,  $D$  denote our data set of  $N$  observations,  $G$  be a DBN, and  $\mathbf{Pa}_i = \{X_{i1}, \dots, X_{is_i}\}$  be the set of parents of  $X_i$  in  $G$  with corresponding  $\{r_{i1}, \dots, r_{is_i}\}$  discrete states,  $s_i = |\mathbf{Pa}_i|$ . The MIT score is then defined as:

$$S_{MIT}(G : D) = \sum_{i=1; \mathbf{Pa}_i \neq \emptyset}^n \{2N \cdot I(X_i, \mathbf{Pa}_i) - \sum_{j=1}^{s_i} \chi_{\alpha, l_{i\sigma_i(j)}}\}$$

where  $I(X_i, \mathbf{Pa}_i)$  is the mutual information between  $X_i$  and its parents as estimated from  $D$ .  $\chi_{\alpha, l_{ij}}$  is the value such that  $p(\chi^2(l_{ij}) \leq \chi_{\alpha, l_{ij}}) = \alpha$  (the Chi-square distribution at significance level  $1 - \alpha$ ), and the term  $l_{i\sigma_i(j)}$

is defined as:

$$l_{i\sigma_i(j)} = \begin{cases} (r_i - 1)(r_{i\sigma_i(j)} - 1) \prod_{k=1}^{j-1} r_{i\sigma_i(k)}, & j = 2 \dots, s_i \\ (r_i - 1)(r_{i\sigma_i(j)} - 1), & j = 1 \end{cases}$$

where  $\sigma_i = \{\sigma_i(1), \dots, \sigma_i(s_i)\}$  is any permutation of the index set  $\{1 \dots s_i\}$  of  $\mathbf{Pa}_i$ , with the first variable having the greatest number of states, the second variable having the second largest number of states, and so on. MIT falls under the same category of information theory based scores as the MDL, Bayesian Information Criterion (BIC) and Akaike Information Criterion (AIC). Briefly speaking, under MIT, the goodness-of-fit of a network is measured by the total mutual information shared between each node and its parents, penalized by a term which quantifies the degree of statistical significance of this shared information. Through extensive experimental validation, de Campos (2006) suggested that, for the task of learning static BN, MIT can compete favorably with Bayesian scores (BDe), outperforms BIC/MDL and should be the score of reference within those based on information theory. However, as opposed to the other popular scoring metrics, to our knowledge MIT has not been considered for DBN learning.

In our recent work (Vinh et al., 2011), we have shown that under the same set of assumptions made in Dojer (2006), there exists a polynomial worst-case time complexity algorithm for learning the globally optimal DBN structure with MIT. We call this algorithm GlobalMIT. The polynomial worst-case time complexity of GlobalMIT is characterized by:

$$p_{MIT}^* = \arg \min \{p \mid \sum_{j=1}^p \chi_{\alpha, l_{i\sigma_i(j)}} \geq 2N \cdot \log k\},$$

It can be seen that  $p_{MIT}^*$  depends only on  $\alpha, k$  and  $N$ . In the worst case, our algorithm will have to examine all the possible parent sets of cardinality from 1 to  $p_{MIT}^* - 1$ . Since there are  $O(n^{p_{MIT}^*})$  subsets with at most  $p_{MIT}^* - 1$  parents, and each set of parents can be scored in polynomial time, globalMIT admits an overall polynomial worst-case time complexity in the number of variables (see also *GlobalMIT user guide* for further details). Our experimental evaluation in Vinh et al. (2011) showed that GlobalMIT is very competitive in terms of network quality. In other words, GlobalMIT seems to combine the strength of both MDL (speed) and BDe (solution quality).

## 2.1 Implementation

The algorithm involves investigating, for each variable, every potential parent set of increasing cardinality (not exceeding  $p_{MIT}^* - 1$ ) until the globally optimal solution has been found. One important observation for the efficient implementation of GlobalMIT is the following decomposition property of the mutual information:

$$I(X_i, \mathbf{Pa}_i \cup X_j) = I(X_i, \mathbf{Pa}_i) + I(X_i, X_j | \mathbf{Pa}_i)$$

This implies that the mutual information can be computed incrementally, and suggests that, for efficiency, the computed mutual information values should be cached to avoid redundant computations.

We provide an implementation of the GlobalMIT algorithm as a Matlab toolbox. The toolbox also supports simple data pre-processing functionalities, such as data discretization and mapping, and simple post-processing such as visualization and quality assessment. For improved performance, we also provide an implementation of GlobalMIT in C++. Our experiments showed that the C++ version of GlobalMIT is up to 40 times on average faster than the Matlab version. For seamless and easy use of GlobalMIT, interface modules provide connection between Matlab and the C++ version, allowing the users to perform all operations in Matlab.

We experimentally compared the runtime of GlobalMIT to BNFinder (Wilczynski and Dojer, 2009) with both the MDL and BDe metrics. The test was carried out on a synthetic data set of 20 genes  $\times$  2000 observations, generated from a gene regulatory network (Network No. 1) as described in (Yu et al., 2004). On a Core 2 Duo PC with 4Gb of main memory, GlobalMIT C++ and BNFinder+MDL took slightly more than an hour to analyze this data set, while BNFinder+BDe took more than three days.

It is noted that currently GlobalMIT only learns DBN with inter-time slice edges, i.e., edges from  $X_i^{(t-1)}$  to  $X_j^{(t)}$ . Learning DBN which allows both inter and intra-time slice edges falls back to an NP-hard problem. However, intra-time slice edges representing (almost) instantaneous genetic interactions, if of interest, can be learned separately using some BN learning algorithm, then combine with the inter-time slice edges, followed by some post processing for the final result. Our future work includes expanding GlobalMIT in this direction. Also, we are extending our framework to handle edges spanning either two or several time slices, which represent variable, longer time-delayed genetic interactions that are also abundant in genetic networks (Zou and Conzen, 2005).

## ACKNOWLEDGEMENT

This research forms part of a project supported by an Australia-India strategic research fund (AISRF).

## REFERENCES

- Chickering, D. M. (1996). Learning Bayesian Networks is NP-Complete. In D. Fisher and H. Lenz, editors, *Learning from Data: Artificial Intelligence and Statistics V*, pages 121–130.
- de Campos, L. M. (2006). A scoring function for learning bayesian networks based on mutual information and conditional independence tests. *J. Mach. Learn. Res.*, 7, 2149–2187.
- Dojer, N. (2006). Learning Bayesian Networks Does Not Have to Be NP-Hard. In *Proceedings of International Symposium on Mathematical Foundations of Computer Science*, pages 305–314.
- Husmeier, D. (2003). Sensitivity and specificity of inferring genetic regulatory interactions from microarray experiments with dynamic Bayesian networks. *Bioinformatics*, 19(17), 2271–2282.
- Kim, S. Y., Imoto, S., and Miyano, S. (2003). Inferring gene networks from time series microarray data using dynamic Bayesian networks. *Briefings in Bioinformatics*, 4(3), 228–235.
- Murphy, K. and Mian, S. (1999). Modelling gene expression data using dynamic bayesian networks. Technical report, Computer Science Division, University of California, Berkeley, CA.
- Perrin, B.-E., Ralaivola, L., Mazurie, A., Bottani, S., Mallet, J., and dAlchBuc, F. (2003). Gene networks inference using dynamic Bayesian networks. *Bioinformatics*, 19(suppl 2), ii138–ii148.
- Scutari, M. (2010). Learning Bayesian Networks with the bnlearn R Package. *Journal of Statistical Software*, 35(3), 1–22.
- Smith, V. A., Yu, J., Smulders, T. V., Hartemink, A. J., and Jarvis, E. D. (2006). Computational inference of neural information flow networks. *PLoS Comput Biol*, 2(11), e161.
- Vinh, N. X., Chetty, M., Coppel, R., and Wangikar, P. P. (2011). A polynomial time algorithm for learning globally optimal dynamic Bayesian network and its applications in genetic network reconstruction. Technical report, Faculty of IT, Monash University.
- Wilczynski, B. and Dojer, N. (2009). BNFinder: exact and efficient method for learning Bayesian networks. *Bioinformatics*, 25(2), 286–287.
- Yu, J., Smith, V. A., Wang, P. P., Hartemink, A. J., and Jarvis, E. D. (2004). Advances to Bayesian network inference for generating causal networks from observational biological data. *Bioinformatics*, 20(18), 3594–3603.
- Zou, M. and Conzen, S. D. (2005). A new dynamic Bayesian network (DBN) approach for identifying gene regulatory networks from time course microarray data. *Bioinformatics*, 21(1), 71–79.