# Modelling Non-Stationary Gene Regulatory Process with Hidden Markov Dynamic Bayesian Network

Shijia Zhu

School of Computer Science and Technology
Harbin Institute of Technology,
Harbin, Heilongjiang, China
sjzhu@hit.edu.cn

Yadong Wang

School of Computer Science and Technology
Harbin Institute of Technology,
Harbin, Heilongjiang, China
ydwang@hit.edu.cn

*Abstract*—**Dynamic Bayesian Network (DBN) has been widely used to infer gene regulatory network from time series gene expression dataset. The standard assumption underlying DBN is based on stationarity, however, in many cases, the gene regulatory network topology might evolve over time. In this paper, we propose a novel non-stationary DBN based network inference approach. In this model, for each variable, a specific HMM implicitly well handles the transition of the stationary DBNs along timesteps. Furthermore, we present a criterion, named as BWBIC score. This criterion is an approximation to the EM objective term, which can reasonably and easily evaluate hmDBN. Towards BWBIC score, a greedy hill climbing based structural EM algorithm is proposed to efficiently infer the hmDBN model. We respectively apply our method on synthetic and real biological data. Compared to the recent proposed methods, we obtained better prediction accuracy on both datasets.**

*Keywords-gene regulatory network; hmDBN; non-stationary DBN; HMM; DBN*

## I. INTRODUCTION

The standard assumption underlying traditional DBN is stationarity, i.e. the structure and parameter of DBN are fixed throughout the time. However, this assumption is not true in many cases. Recently, various research efforts have therefore addressed the stationary assumption for DBNs, such as time-varying graphical models with known priori about the number of time segments [1-2], non-stationary graphical models with global transition times [3-4], node-specific undirected networks [5] and continuous non-stationary DBN with globally fixed network topologies [6].

In this paper, we present a node specific non-stationary network by combining HMM with DBN. This model is called *hidden markov Dynamic Bayesian Network* (*hmDBN*). This new framework incorporates the notion of Baum-Welch algorithm. It is suitable for studying the problem, in which the network structure is varying slowly and smoothly over time. On this basis, we present a novel score criterion for evaluating a specific hmDBN (called BWBIC score). This criterion is an approximation to the EM objective term involving the graph structure, which can be calculated easily. Furthermore, we proposed a novel greedy hill climbing based structural EM algorithm to recover the hmDBN model. It is an iterative approach (called hmDBN structure learning algorithm), which selects appropriate hmDBN structure and

transition matrix parameters to maximize the objective term given observations.

## II. METHODS

### A. *The hmDBN Model*

The *hmDBN* is obtained by combining HMM with DBN. It employs the transitions among hidden states of HMM to describe the evolution of network structure over time, and utilizes DBN to describe the conditional dependence among variables. Based on conditional independence for DBN, the *hmDBN* can be decomposed into non-stationary sub-networks for all variables $X=\{X_1, X_2, ..., X_N\}$, i.e. $hmDBN=\{hmDBN_1, hmDBN_2, ..., hmDBN_N\}$. In $hmDBN_i$, only the structure among variable $X_i$ and its parents is retained from non-stationary DBN. Thus, it allows different transition times for the non-stationary sub-networks of different variables.
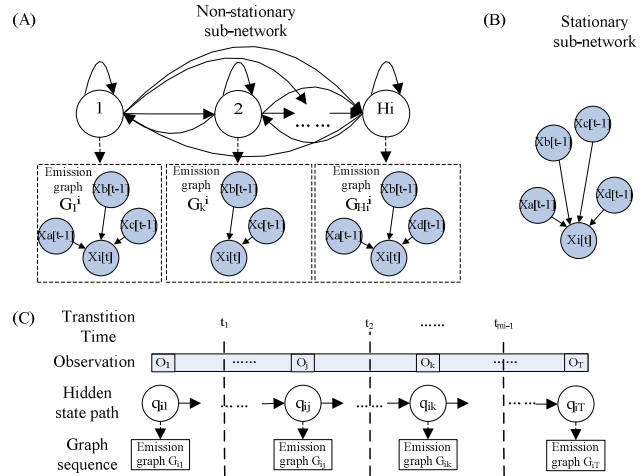


Figure 1. Graphical structure of hmDBN for a specific variable $X_i$. (A) Shaded nodes represent observed nodes. The blank nodes denote the hidden states, which could transmit among different states over time. One hidden state corresponds to one network topology in the rectangle (emission graph). (B) The stationary sub-network for variable corresponding to the non-stationary sub-network for variable $X_i$. (C) The hidden state path indicates the network structure at each timestep.

The evolution of the non-stationary sub-network for each specific variable follows a Markov assumption. Fig. 1A is a graphical illustration of non-stationary sub-network for a specific variable $X_i$ (*hmDBN_i*). The fully connected blank

nodes represent the hidden states. Each hidden state corresponds to one graph in the rectangle (called an *emission graph*). Each emission graph is a DBN, in which only $X_i$ is child node. The number of hidden states is identical to the number of emission graphs. In addition, we define the network which consists of all nodes and edges appearing in all emission graphs of $X_i$ as the *stationary sub-network* of $X_i$. The graph in Fig. 1B shows the stationary sub-network corresponding to the non-stationary sub-network in Fig. 1A. Given an observation sequence $O=\{O_1,O_2,...,O_T\}$, its corresponding hidden state path reflects the emission graph sequences of the non-stationary sub-network for variable $X_i$, i.e. $G^i=\{G_1^i,\ G_2^i,...,G_T^i\}$. Herein, we use notation $q^i=\{q_1^i,\ q_2^i,...,q_T^i\}$ to represent the hidden state path (Fig. 1C). $hmDBN_i$ for variable $X_i$ takes the following form:

$$P\left(G^i,O\mid hmDBN_i\right)=P\left(G^i\mid hmDBN_i\right)P\left(O\mid G^i,hmDBN_i\right)$$

$$=\pi_{q_1}^i\left(\prod_{t=1}^T p\left(O_t\mid G_{q_t}^i\right)\right)\left(\prod_{t=2}^T a_{(q_{t-1})q_t}^i\right) \tag{1}$$

where $\pi_{q_1}^i=P\left(G_{q_1}^i\right)$ represents the distribution for the initial state $q_1^i$, and $A^i=\{a_{q_a,q_b}^i\}$ represents the transition matrix.

## B.  The BWBIC Score

We utilized Expectation-Maximization (EM) algorithm to estimate the parameters for $hmDBN_i$. EM algorithm searches for the value of the following function (For brevity, we omit the index $i$ of parameters for $hmDBN_i$):

$$Q(\lambda,\lambda')=E\left[\log p\left(X_{obs},X_{mis}\mid\lambda\right)\mid X_{obs},\lambda'\right]$$

$$\propto\sum_q\log\pi_{q_1}P(O,q\mid\lambda')+\sum_q\left(\sum_{t=2}^T\log a_{(q_{t-1})q_t}\right)P(O,q\mid\lambda')+$$

$$\sum_q\left(\sum_{t=1}^T\log p\left(O_t\mid G_{qt}\right)\right)P(O,q\mid\lambda') \tag{2}$$

where observable value $X_{obs}$ is temporal sequence data $O=\{O_1,...,O_T\}$; missing value $X_{mis}$ is hidden state path $q$; $\lambda=\{\pi,A,G\}$ and $\lambda'$ are the current estimates of $\lambda$.

The first two terms in formula (2) can be maximized by optimizing $\lambda'=\arg\max_\lambda Q(\lambda,\lambda')$. The parameters $\pi$ and $A$ can be re-estimated as following:

$$\pi_j'=p\left(O,q_1=j\mid\lambda'\right)\Big/\sum_{j=1}^H p\left(O,q_1=j\mid\lambda'\right) \tag{3}$$

$$a_{jk}'=\sum_{t=2}^T p\left(O,q_{t-1}=j,q_t=k\mid\lambda'\right)\Big/\sum_{t=2}^T p\left(O,q_{t-1}=j\mid\lambda'\right) \tag{4}$$

However, the third term is very hard to maximize by solving out the parameter $G$, since it is a NP-complete problem. Indeed, the third term is still difficult to calculate, given a network structure. Therefore, we give an easily calculated asymptotic approximation of the third term as following:

$$\sum_q\left(\sum_{t=1}^T\log p\left(O_t\mid G_{qt}\right)\right)P(O,q\mid\lambda')$$

$$\approx\sum_{g=1}^H\left(\log F\left(O\mid G_g,\theta_g^*\right)-\frac{d_g}{2}\log m_g\right) \tag{5}$$

$$\log F\left(O\mid G_g,\theta_g^*\right)$$

$$=\sum_{j=1}^{S_g}\sum_{k=1}^R\sum_{t=1}^T P(O,q_t=g\mid\lambda')\chi\left(i,j,k:O_t\right)\log\theta_{g,ijk}^* \tag{6}$$

$$m_g=\sum_{t=1}^T P(O,q_t=g\mid\lambda') \tag{7}$$

$$\theta_{g,ijk}=P\left(X_i=k\mid\pi\left(X_i\right)=j,G_g\right) \tag{8}$$

$$\theta_{g,ijk}^*=\frac{\sum_{t=1}^T\chi\left(i,j,k:O_t\right)P(O,q_t=g\mid\lambda')}{\sum_{k=1}^R\sum_{t=1}^T\chi\left(i,j,k:O_t\right)P(O,q_t=g\mid\lambda')} \tag{9}$$

$$\chi\left(i,j,k:O_t\right)=\begin{cases}1,if\ \ X_i=k\ \ and\ \ \pi\left(X_i\right)=j,in\ \ O_t\\0,otherwise\end{cases} \tag{10}$$

where $H$ is the number of the hidden states (emission graphs) of $hmDBN_i$; $R$ is the number of values of variable $X_i$; $S_g$ is the number of values of the parent nodes $\pi\left(X_i\right)$ in $g$-th emission graph with $1\le g\le H$; $G_g$ represents the $g$-th emission graph in $hmDBN_i$; $\theta_g$ represent the network parameters of emission graph $G_g$; $\theta_g$ are independent and multinomially distributed; they are composed of $d_g$ independent parameters, and $d_g=S_g(R-1)$. $P(O,q_t=g\mid\lambda')$ represents the probability that the observation values at timestep $t$ are generated by the $g$-th hidden state, and we call it *sample distribution* for observation $O$. It can be calculated using forward and backward algorithms [7].

## C.  The hmDBN Structure Learning Algorithm

We proposed an algorithm to maximize the approximated Q function based on BWBIC score by selecting appropriate transition matrix parameters and non-stationary structure. We call it *hmDBN structure learning algorithm*. The steps of the algorithm are as following:

**1)** On the basis of former stationary network, randomly generate a new stationary network.

**2)** transform the stationary DBN into hmDBN;

**2.1)** set initial values for hmDBN, including putative emission graph structures, transition matrix parameters and initial hidden state path;

**2.2)** based on the initial values, re-estimate the transition matrix parameters according to the formula (3-4);

**2.3)** estimate the hidden state path using Viterbi algorithm [7];

**3)** calculate the BWBIC score according to formula (5);

**4)** use a greedy hill climbing strategy to assess whether the BWBIC score is better than the score for previous structure, and determines whether the search converges. If it does not converge, step 1) generates another stationary network, and the above steps repeat.

These steps comprise the EM algorithm. Step 1), 2.1), 2.2) and 4) comprise the M-step, which respectively updates the transition matrix and selects the appropriate non-stationary DBN structure, thereby improving the Q function. Step 2.3) and 3) comprise the E-step, which respectively calculates the Q function and estimates the values of hidden variables. This algorithm does not simultaneously optimize the network structure and parameters, but rather optimizes the parameters given the fixed network structure, and next, optimizes the network structure. Either of the two steps can improve the Q function. Our proposed algorithm is a structural EM approach [8]. Our contribution is that we construct the BWBIC-based structural EM Q function for the hmDBN model and design the corresponding maximization method.
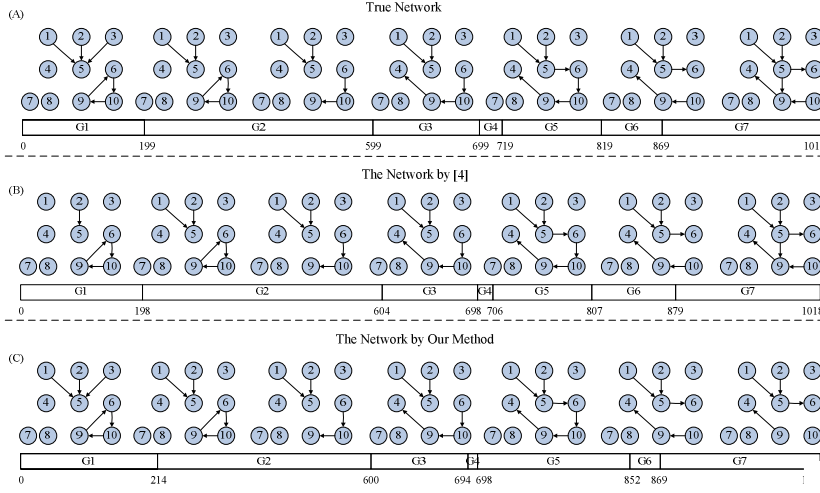
Figure 2. Structure learning for a simulated data: (A) True non-stationary DBN. (B) The non-stationary DBN by [4], under the settings of known transition number and unknown transition times. (C) The non-stationary DBN constructed by our proposed method without knowing about transition times and epoch number.



Figure 4. Sample distribution for the genes whose regulators change over time: Horizontal axis denotes the time steps. Vertical axis represents the probability density. The curves in red and blue respectively represents the sample distributions over two graphs in the same color on the left of vertical axis. The green dash line indicates the transition times.
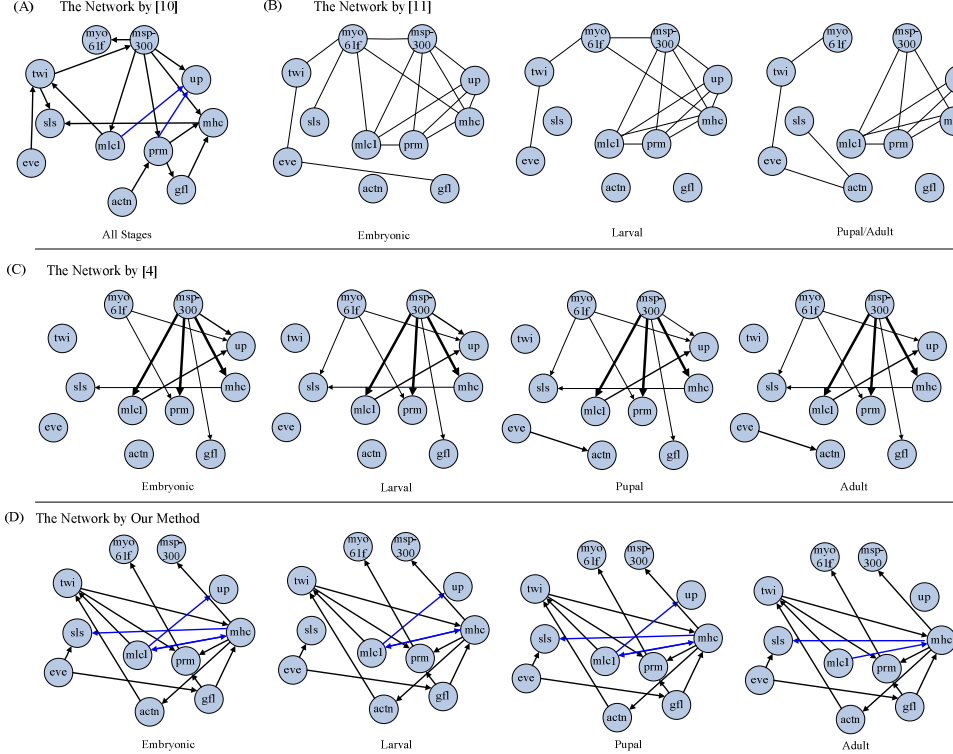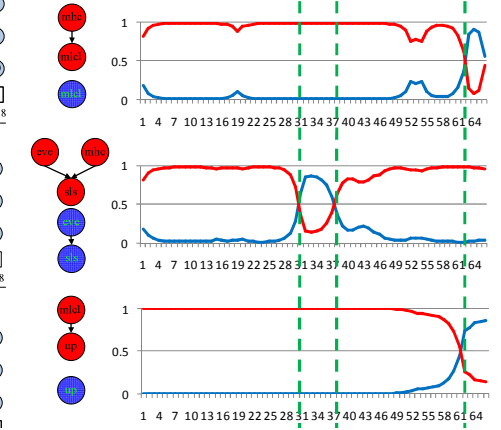


Figure 3. Structure learning for Drosophila muscle development data: (A) The network for all stages reported by [10]. (B) The undirected network reported by [11]. (C) The network reported by [4] under the settings of known epoch number and known transition times. (D) The non-stationary DBN constructed using our proposed method. The edges in blue are the different edges among different segments.

## III. RESULTS

### A. *Evaluation Using Simulated Data*

We performed the proposed approach on a dataset generated by the same method as in [4]. This dataset is on a ten node network with six single-edge changes between seven epochs where the length of each segment varies between 20 and 400 observations (Fig. 2A). Under the settings of unknown transition times and unknown transition number, the study [4] predicted six epochs with the maximum probability about 0.5. Our algorithm successfully recovers seven segments. On the other hand, the network structures reconstructed by [4] and us are

shown in Fig. 2B and 2C respectively. Both two structures are very close to the actual network topology. However, it is worth noting that the structure recovered by [4] is obtained under the setting of known transition number, while our network structure is identified without knowing the transition number. Thus, our algorithm obtains better results, even if we have less knowledge in advance.

## B. Evaluation Using Gene Expression Data

We also apply our method to recover non-stationary networks using the Drosophila development gene expression data from [9]. This data contains expression measurements over 66 time steps of 4028 Drosophila genes throughout development and growth during the embryonic, larval, pupal, and adult stages of life. We choose 11 genes for analysis, which are *eve*, *gfl/lmd*, *twi*, *mlc1*, *sls*, *mhc*, *prm*, *actn*, *up*, *myo61f* and *msp300*. Those genes were reported to be related with the muscle development of Drosophila. We preprocessed continuous expression data into binary values using the same methods as in [10]. Our algorithm takes about 2 minutes for this dataset on a 3.2GHz Intel Core i5 machine with 4 GB of RAM, and this procedure is executed without parallel.

Using these genes, some other researchers have employed different approaches to reconstruct the regulatory networks, such as a stationary directed network reconstructed by [10], a non-stationary undirected network predicted by [11] and a non-stationary DBN inferred by [4] (Fig. 3). In those predictions, a cluster forms around *myo61f*, *msp-300*, *up*, *mhc*, *prm*, and *mlc1*. All of these genes except *up* are in the myosin family, which contains genes involved in muscle contraction. Our prediction also indicates the intense associations among these genes, although the inferred relationships are different from other predictions. Additionally, we demonstrate that the *up* gene is not so intensely connected with this cluster as other genes, and *up* gene disassociated with this cluster in the adult stage. Currently, the reference regulatory network on the muscle development of Drosophila is not available and the relevant biological literatures are limited. The only existing evidence includes two aspects. First, the transition times of four life periods are located at 30, 40 and 58. The study [11] predicts 3 segments, and Robinson et al. predict 4 segments with the posterior peaks located at 30, 40 and 50. Our prediction indicated four segments with the transition time 30, 37 and 61, which is even closer to the experiments. Fig. 4 presents the sample distribution for the gene whose regulators change over the time. The green dash lines indicate the transition times for the non-stationary network. Second, *gfl/lmd* and *twi* are indicated to direct upstream regulators of mef2 [12-13] that directly regulates some target myosin family genes at all stages of muscle development [14], such as *mhc* and *mlc1*. The study [11] did not predict the interaction of gene *twi* and *gfl* with gene *mhc* and *mlc1*, and under the settings of known transition time and transition number, Robinson et al. [4] still miss such association. However, our algorithm predicted these interactions, even if we do not know the transition time and transition number.

## IV.    CONCLUSION

In this paper, we present a node specific non-stationary network (hmDBN) that combines two efficient and well-tried techniques: HMM with DBN. This new framework is suitable for studying the problem, in which the network structure is varying slowly and smoothly over time. On this basis, we present a novel score criterion for evaluating a specific hmDBN (called BWBIC score). Furthermore, we proposed a novel greedy hill climbing based structural EM algorithm to recover the hmDBN model. We apply our method to two datasets. Compared to recent proposed methods, we obtain better performance in both cases.

## REFERENCES

[1]      R. Yoshida, S. Imoto, and T. Higuchi, "Estimating time-dependent gene networks from time series microarray data by dynamic linear models with Markov switching," *Proc IEEE Comput Syst Bioinform Conf*, pp. 289-98, 2005.

[2]      M. Talih, and N. Hengartner, "Structural learning with time-varying components: tracking the cross-section of financial time series," *Journal of the Royal Statistical Society Series B,* vol. 67, no. 3, pp. 321-341, 2005.

[3]      X. Xuan, and K. P. Murphy, "Modeling changing dependency structure in multivariate time series," in Proceedings of the 24th International Conference on Machine Learning (ICML 2007), 2007, pp. 1055-1062.

[4]      J. Robinson, and A. Hartemink, "Non-Stationary Dynamic Bayesian Networks," in Neural Information Processing Systems 2008 (NIPS 2008), 2008, pp. 1369-1376.

[5]      A. Ahmed, and E. P. Xing, "Recovering time-varying networks of dependencies in social and biological studies," *Proc Natl Acad Sci U S A,* vol. 106, no. 29, pp. 11878-83, Jul 21, 2009.

[6]      M. Grzegorczyk, and D. Husmeier, "Non-stationary continuous dynamic Bayesian networks," in Neural Information Processing Systems 2009 (NIPS 2009), 2009, pp. 682-690.

[7]      L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition " *Proceedings of the IEEE,* vol. 77, no. 2, pp. 257 - 286 1989.

[8]      N. Friedman, "The Bayesian structural EM algorithm," in UAI'98 Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence, 1998, pp. 129-138.

[9]      M. N. Arbeitman, E. E. Furlong, F. Imam *et al.*, "Gene expression during the life cycle of Drosophila melanogaster," *Science,* vol. 297, no. 5590, pp. 2270-5, Sep 27, 2002.

[10]      W. Zhao, E. Serpedin, and E. R. Dougherty, "Inferring gene regulatory networks from time series data using the minimum description length principle," *Bioinformatics,* vol. 22, no. 17, pp. 2129-35, Sep 1, 2006.

[11]      F. Guo, S. Hanneke, W. Fu *et al.*, "Recovering temporally rewiring networks: A model-based approach," in Proceedings of the 24th International Conference on Machine Learning (ICML 2007), 2007, pp. 321-328.

[12]      R. M. Cripps, B. L. Black, B. Zhao *et al.*, "The myogenic regulatory gene Mef2 is a direct target for transcriptional activation by Twist during Drosophila myogenesis," *Genes Dev,* vol. 12, no. 3, pp. 422-34, Feb 1, 1998.

[13]      H. Duan, and H. T. Nguyen, "Distinct posttranscriptional mechanisms regulate the activity of the Zn finger transcription factor lame duck during Drosophila myogenesis," *Mol Cell Biol,* vol. 26, no. 4, pp. 1414-23, Feb, 2006.

[14]      T. Sandmann, L. J. Jensen, J. S. Jakobsen *et al.*, "A temporal map of transcription factor activity: mef2 directly regulates target genes at all stages of muscle development," *Dev Cell,* vol. 10, no. 6, pp. 797-807, Jun, 2006.