

Learning a Probabilistic Boolean Network Model from Biological Pathways and Time-series Expression Data

A thesis submitted in partial fulfillment of the requirements for the degree of

Bachelor of Technology

in

Electronics and Electrical Communication Engineering

by

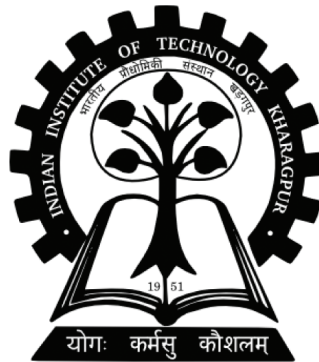
Vardaan Pahuja (12EC10067)

advised by

Prof. Ritwik Kumar Layek

and

Prof. Pabitra Mitra



Department of Electronics and Electrical Communication Engineering
Indian Institute of Technology, Kharagpur

April 2016

Certificate

This is to certify that the thesis titled **Learning a Probabilistic Boolean Network Model from Biological Pathways and Time-series Expression Data** submitted by **Vardaan Pahuja (12EC10067)** to the Department of Electronics and Electrical Communication Engineering is a bonafide record of work carried out by him under my supervision and guidance. The thesis has fulfilled all the requirements as per the regulations of the Institute and, in my opinion, has reached the standard needed for submission.

Prof. Ritwik Kumar Layek

Assistant Professor

Department of Electronics and

Electrical Communication Engineering

IIT Kharagpur

Prof. Pabitra Mitra

Associate Professor

Department of Computer Science

and Engineering

IIT Kharagpur

Declaration

I, **Vardaan Pahuja** hereby certify that this report titled **Learning a Probabilistic Boolean Network Model from Biological Pathways and Time-series Expression Data**

- Is an original work and has been done by me under the guidance of my supervisor;
- The work has not been submitted to any other Institute for any degree or diploma;
- While writing the report I have conformed to norms and guidelines given in the Ethical Code of Conduct of the Institute;
- Whenever I have used materials (data, model, figures and text) from other sources, I have given due credit to them by citing them in the text of the report, giving their details in the references, and following fair use doctrine policies of copy righted materials if any used in this thesis.

Vardaan Pahuja (12EC10067)

12th April, 2016

Acknowledgement

I sincerely thank **Prof. Ritwik Kumar Layek**, Department of Electronics and Electrical Communication Engineering and **Prof. Pabitra Mitra**, Department of Computer Science and Engineering, IIT Kharagpur for continuously supporting me throughout the duration of the project. Their valuable inputs were of immense help to solve the challenges involved in the project. They have taken pain to go through my work and make necessary corrections as and when needed. Their consistent encouragement and co-operation during the project is really invaluable.

Vardaan Pahuja (12EC10067)

12th April, 2016

Abstract

The problem of inferring a stochastic model for gene regulatory networks is addressed here. The prior biological data includes biological pathways and time-series expression data. We propose a novel algorithm to use both of these data to construct a Probabilistic Boolean Network (PBN) which models the observed dynamics of genes with a high degree of precision. Our algorithm constructs a pathway tree and uses the time-series expression data to select an optimal level of tree, whose nodes are used to infer the PBN.

Contents

Contents	i
1 Introduction	1
2 Preliminaries	2
2.1 Boolean Networks	2
2.1.1 Introduction	2
2.1.2 Restricted Boolean Networks	2
2.2 Probabilistic Boolean Networks	3
2.2.1 Introduction	3
2.2.2 Gene Influence in PBN	4
2.3 Inferring Regulatory relationships using Biological Pathways	4
2.3.1 Introduction	4
2.3.2 Algorithm	5
2.3.3 Reduction of Boolean search space	5
2.4 Inferring Regulatory relationships using time-series expression data	6
2.4.1 Three-rule method	6
2.4.2 Constraint based analysis of regulatory relationships	7
2.4.3 Error Analysis	9
2.4.4 Inference Algorithm	9
3 Proposed Algorithm	11
3.1 Introduction	11
3.2 Construction of Pathway Tree	11

<i>CONTENTS</i>	iii
3.3 Selecting the optimum level of tree	12
4 Performance Evaluation	14
5 Results and Discussion	16
6 Inference of Yeast Cell Cycle Network	19
7 Conclusion	22

Chapter 1

Introduction

Modelling cellular interaction dynamics has been one of the important issues in systems biology [2]. A number of mathematical formulations have been proposed to model these genetic interactions, including Bayesian networks [6], linear models [10], and Boolean networks [3]. Based on a couple of limitations of BNs (e.g. limitation that BN is a deterministic model), a stochastic version of BNs, i.e. Probabilistic Boolean Networks (PBN) was proposed by Shmulevish et al. [9]. It incorporates uncertainty both in data and model selection.

The task of inferring gene regulatory networks from prior biological data is an ill-posed inverse problem, since multiple network realizations could explain the same biological phenomenon. The search space for potential regulatory genes and the boolean functions associated with them, increases exponentially with the number of genes in the network. Use of biological pathways to infer boolean networks was demonstrated in [4]. Restricted boolean networks are simplified boolean networks in which the regulatory relationships between genes is either activation (positive regulation to target gene) or inhibition (negative regulation to target gene). A three-rule method to construct a restricted Boolean network from time-series data was proposed by Higa et al. [1].

Here, we propose a novel algorithm which utilizes both biological pathway data and time-series expression data to construct a Probabilistic Boolean network to model the gene dynamics. Earlier approaches use a single form of biological data, which could be subjected to experimental bias. We overcome this limitation in our algorithm by using two different forms of data to infer the PBN.

Chapter 2

Preliminaries

2.1 Boolean Networks

2.1.1 Introduction

A Boolean network (BN) $G(V, F)$ on n genes is defined by a set of nodes/genes such that each node has a Boolean function assigned to it. Here F is the set of Boolean functions where,

$$F = \{f_1, f_2, \dots, f_n\}, f_i : \{0, 1\}^n \rightarrow \{0, 1\}, \quad (2.1)$$

and V is the set of nodes, $V = \{v_1, v_2, \dots, v_n\}$. The value v_i denotes the state of gene i , which can be either 0(off) or 1(on). The dynamics of BN can be expressed as,

$$v_i(t+1) = f_i(v_1(t), v_2(t), \dots, v_n(t)) = f_i(\mathbf{v}(t)) \quad (2.2)$$

Here $\mathbf{v}(t)$ is called the Gene Activity Profile (GAP).

2.1.2 Restricted Boolean Networks

Restricted Boolean networks are simplified Boolean networks in which the regulatory relationships between genes obey the following convention: $a_{ij} = 1$ represents a positive regulation from gene x_j to x_i (activation); $a_{ij} = -1$ represents a negative regulation from gene

x_j to x_i (inhibition); and $a_{ij} = 0$ means that x_j has no effect on x_i . The Boolean function $f_i(x_1, \dots, x_{k_i})$ is defined as [5]

$$x_i(t+1) = \begin{cases} 1, & \text{if } \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_j(t) > 0 \\ 0, & \text{if } \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_j(t) < 0 \\ x_i(t), & \text{if } \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_j(t) = 0 \end{cases} \quad (2.3)$$

2.2 Probabilistic Boolean Networks

2.2.1 Introduction

BN is a deterministic model. However, due to inherent uncertainty associated with a biological system, a stochastic model is more appropriate here [9]. Probabilistic Boolean Network is a stochastic version of BN in which more than one Boolean function can be assigned to a gene. Thus, for every node, there corresponds a set

$$F_i = \{f_j^{(i)}\}_{j=1,2,\dots,l(i)} \quad (2.4)$$

where each $f_j^{(i)}$ is a possible predictor function for gene i and $l(i)$ is the number of possible functions for gene i . The probability of choosing the j^{th} predictor function for gene i is c_i^j . This implies that

$$\sum_{j=1}^{l_i} c_i^j = 1, 0 < c_i^j < 1, \text{ for } i = 1, 2, \dots, n \quad (2.5)$$

If we choose the j_i^{th} Boolean function for gene v_i , then the BN can be expressed as $BN_{j_1, j_2, \dots, j_n}$ where $j_i \in \{1, 2, \dots, l_i\}$. The probability of choosing $BN_{j_1, j_2, \dots, j_n}$ is given by

$$P\{f_1 = f_1^{j_1}, f_2 = f_2^{j_2}, \dots, f_n = f_n^{j_n}\} = \prod_{i=1}^n c_i^{j_i} = q_{j_1 j_2 \dots j_n} \quad (2.6)$$

2.2.2 Gene Influence in PBN

Different genes can have a varying degree of impact on the predictor function of a gene. The partial derivative of a Boolean function with respect to variable x_j ($1 \leq j \leq n$) is defined as

$$\frac{\partial f}{\partial x_j} = f(x^{(j,0)}) \oplus f(x^{(j,1)}) \quad (2.7)$$

where \oplus is modulo-2 addition operation.

The influence of the variable x_i on function f_i is the expectation of the partial derivative with respect to initial joint probability distribution $D(x)$, $x \in \{0, 1\}^n$.

$$\begin{aligned} I_j(f) &= E_D \left[\frac{\partial f}{\partial x_j} \right] = \Pr \left\{ \frac{\partial f}{\partial x_j} = 1 \right\} \\ &= \Pr [f(x) \neq f(x^{(j)})] \end{aligned} \quad (2.8)$$

where $x^{(j)}$ is same as x except that the j^{th} component is toggled.

2.3 Inferring Regulatory relationships using Biological Pathways

2.3.1 Introduction

The pathway segment $A \xrightarrow{t:a,b} B$ implies that if gene A assumes the value a , then gene B transitions to b in no more than t subsequent time-stamps [4]. A pathway is defined to be a sequence of pathway segments of the form $A \xrightarrow{t_1:a,b} B \xrightarrow{t_2:b,c} C$. A trajectory is a sequence of states $S_0 \rightarrow S_1 \rightarrow S_2 \rightarrow S_3 \rightarrow S_4$ resulting from network rules beginning at some initial state. These pathways represent *a priori* biological information. Our goal is to generate a PBN whose trajectories are consistent with the given set of biological pathways. This is an ill-posed inverse problem that could have multiple solutions or perhaps none.

2.3.2 Algorithm

1. Consider the Karnaugh maps for each gene state at next time stamp. Initially, this corresponds to the entire space of Boolean networks for n genes.
2. For each pathway, modify the K-map for the output gene to satisfy the constraint pathways. If no conflict with existing K-map arises, repeat the process for next pathway.
3. In case of conflict, fill the non-conflicting minterms in the K-map with the output gene and introduce new pathways such that the conflicting minterms transition to the non-conflicting minterms in the next time-stamp. The new pathways are added to the queue of existing pathways.
4. At some point, if it's not possible to satisfy the pathways or we return to the original conflict, then we terminate by concluding it's not possible to satisfy all constraints using the Boolean network.

2.3.3 Reduction of Boolean search space

We can impose the additional constraint that the maximum number of predictors allowed for each gene is 2. Such an upper limit on the number of predictors per gene could be motivated from the biological consideration that the promoter region for a gene only has enough room for at most only a few transcription factors to bind. Further reduction in the cardinality of the family of networks can be achieved by imposing additional constraints such as the number and relative significance of the attractors, upper bounds on network connectivity, etc.

Table 2.1: Regulatory relationships for one input gene

ID	$x_{j_1}(t)$	$x_i(t) \rightarrow x_i(t+1)$	a_{ij_1}
1	1	$0 \rightarrow 0$	-1
2	1	$0 \rightarrow 1$	1
3	1	$1 \rightarrow 0$	-1
4	1	$1 \rightarrow 1$	1

2.4 Inferring Regulatory relationships using time-series expression data

2.4.1 Three-rule method

A time-series observation can be treated as a trajectory (or random walk) of the state space of the network used to model a real biological system. The three-rule method proposed by Higa et al. [1] is to induce the constraints between genes from the small difference between two similar states and the difference between their next states. Given an m -point time series $S = \{S(1), S(2), \dots, S(m)\}$ of gene expression profiles, where $S(t) \in \{0, 1\}^n$ for $t = 1, 2, \dots, m$, the three rules are as follows:

Rule 1: Let $S(t-1)$, $S(t)$, and $S(t+1)$ be three consecutive states. If $S(t-1)$ and $S(t)$ differ by a single gene x_k , then for each gene x_i such that $x_i(t) \neq x_i(t+1)$, we have x_k directly regulates x_i ; that is, $a_{ik} \neq 0$.

Rule 2: Only the active genes at time t can possibly regulate genes at time $t+1$.

Rule 3: Given two similar states $S(t_1)$ and $S(t_2)$, the difference between $S(t_1+1)$ and $S(t_2+1)$ must result from the genes in their predecessors $S(t_1)$ and $S(t_2)$ that are expressed differently.

Rules 1 and 3 are also applicable to situations where $S(t-1)$ and $S(t)$ or $S(t_1)$ and $S(t_2)$ differ in more than one gene. Cyclically applying these rules to any two states may lead to a group of constraint inequalities between variables a_{ij} .

Table 2.2: Regulatory relationships for two input genes

ID	$x_{j_1}(t)$	$x_{j_2}(t)$	$x_i(t) \rightarrow x_i(t+1)$	a_{ij_1}	a_{ij_2}	Constraint
1	0	1	$0 \rightarrow 0$	No	-1	
2	1	0		-1	No	
3	1	1		-1 or 1	-1 or 1	$a_{ij_1} + a_{ij_2} \leq 0$
4	0	1	$0 \rightarrow 1$	No	1	
5	1	0		1	No	
6	1	1		1	1	
7	0	1	$1 \rightarrow 0$	No	-1	
8	1	0		-1	No	
9	1	1		-1	-1	
10	0	1	$1 \rightarrow 1$	No	1	
11	1	0		1	No	
12	1	1		-1 or 1	-1 or 1	$a_{ij_1} + a_{ij_2} \geq 0$

No: Undetermined ; -1 or 1: Semi-determined

2.4.2 Constraint based analysis of regulatory relationships

Here, we analyze the constraint inequalities in equation (2.3) and use it to infer the regulatory relationships. The target gene can switch its state in four different combinations i.e. $0 \rightarrow 0$, $0 \rightarrow 1$, $1 \rightarrow 0$, and $1 \rightarrow 1$. Only the input genes which are active at time $(t-1)$, contribute to the change of state at time t . Using equation (2.3), the following inequalities are true for different cases:

$$\begin{aligned}
0 \rightarrow 0 : & \quad \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_i(t) \leq 0 \\
0 \rightarrow 1 : & \quad \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_i(t) > 0 \\
1 \rightarrow 0 : & \quad \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_i(t) < 0 \\
1 \rightarrow 1 : & \quad \sum_{j \in \{1, \dots, k_i\}} a_{ij} x_i(t) \geq 0
\end{aligned} \tag{2.9}$$

Table 2.3: Regulatory relationships for three input genes

ID	$x_{j_1}(t)$	$x_{j_2}(t)$	$x_{j_2}(t)$	$x_i(t) \rightarrow x_i(t+1)$	a_{ij_1}	a_{ij_1}	a_{ij_1}	Constraint
1	0	0	1	$0 \rightarrow 0$	No	No	-1	
2	0	1	0		No	-1	No	
3	1	0	0		-1	No	No	
4	0	1	1		No	-1 or 1	-1 or 1	$a_{ij_2} + a_{ij_3} \leq 0$
5	1	0	1		-1 or 1	No	-1 or 1	$a_{ij_1} + a_{ij_3} \leq 0$
6	1	1	0		-1 or 1	-1 or 1	No	$a_{ij_1} + a_{ij_2} \leq 0$
7	1	1	1		-1 or 1	-1 or 1	-1 or 1	$a_{ij_1} + a_{ij_2} + a_{ij_3} < 0$
8	0	0	1	$0 \rightarrow 1$	No	No	1	
9	0	1	0		No	1	No	
10	1	0	0		1	No	No	
11	0	1	1		No	1	1	
12	1	0	1		1	No	1	
13	1	1	0		1	1	No	
14	1	1	1		-1 or 1	-1 or 1	-1 or 1	$a_{ij_1} + a_{ij_2} + a_{ij_3} > 0$
15	0	0	1	$1 \rightarrow 0$	No	No	-1	
16	0	1	0		No	-1	No	
17	1	0	0		-1	No	No	
18	0	1	1		No	-1	-1	
19	1	0	1		-1	No	-1	
20	1	1	0		-1	-1	No	
21	1	1	1		-1 or 1	-1 or 1	-1 or 1	$a_{ij_1} + a_{ij_2} + a_{ij_3} < 0$
22	0	0	1	$1 \rightarrow 1$	No	No	1	
23	0	1	0		No	1	No	
24	1	0	0		1	No	No	
25	0	1	1		No	-1 or 1	-1 or 1	$a_{ij_2} + a_{ij_3} \geq 0$
26	1	0	1		-1 or 1	No	-1 or 1	$a_{ij_1} + a_{ij_3} \geq 0$
27	1	1	0		-1 or 1	-1 or 1	No	$a_{ij_1} + a_{ij_2} \geq 0$
28	1	1	1		-1 or 1	-1 or 1	-1 or 1	$a_{ij_1} + a_{ij_2} + a_{ij_3} > 0$

No: Undetermined ; -1 or 1: Semi-determined

For a single regulatory gene x_{j_1} , these inequalities simplify to $a_{ij_1} = -1$, $a_{ij_1} = 1$, $a_{ij_1} = -1$, and $a_{ij_1} = 1$ respectively. These are presented in Table 2.1. For the case of two regulatory genes, if a single gene is active, then the regulation of the active gene can be inferred but that of the other gene is undetermined. When both of input genes are active, the regulation of both these genes can be determined if the target gene switches its state. In the other case, the relationship between their regulatory relationships is semi-determined because it is governed by a constraint equation. The different cases for two input gene regulatory relationships are presented in Table 2.2. Similar constraint inequalities can be derived for three input gene regulatory relationships as shown in Table 2.3.

Each time series sample gives rise to one of the cases mentioned in the respected tables. Let N_{ij}^{-1} , N_{ij}^1 , and $N_{ij}^{-1,1}$ denote the number of $a_{ij} = -1$, $a_{ij} = 1$, and $a_{ij} = -1$ or 1 respectively. The degree of determination of a regulatory relationship a_{ij} is defined as

$$d_{ij} = |N_{ij}^{-1} - N_{ij}^1| \quad (2.10)$$

Among multiple input genes in a regulatory relationship, the one with the highest d_{ij} is the first to be decided using majority rule. This value is then put into constraint inequalities for inferring other semi-determined relationships. This procedure is then repeated to determine all other regulatory relationships.

2.4.3 Error Analysis

The error arising out of ambiguity in determination of a_{ij} is defined as $\varepsilon_{ij}^{-1,1} = \min(N_{ij}^{-1}, N_{ij}^1)$. Also, the target gene can't switch its state under null input conditions. This error is denoted by ε_i^{null} . The total error of a predictor set is defined as

$$\varepsilon = \varepsilon_i^{null} + \sum_j \varepsilon_{ij}^{-1,1} \quad (2.11)$$

2.4.4 Inference Algorithm

The algorithm used for determining regulatory relationships [8] is given below:

1. Calculate the total error of each combination of one, two, or three regulatory gene sets.
2. Sort the predictor sets in ascending order of their errors.
3. If a gene appears in the first l sets with a frequency greater than or equal to 50%, then it is selected as a regulatory gene.

Chapter 3

Proposed Algorithm

3.1 Introduction

The list of biological pathways satisfied by the biological system is available to us. Further, we assume a priority ordering of these pathways in order of decreasing reliability i.e. pathways higher in order are more accurate than those lower in order. The complete set of pathways can't represent a Boolean network as many these pathways may conflict with each other regarding prediction of a gene output. On the other hand, a subset of pathways represents a family of Boolean networks since the Karnaugh maps representing the BNs can contain several don't care terms. Our proposed algorithm constructs a m-ary tree with each node containing a subset of non-conflicting pathways. The invariant followed is that each child node satisfies the pathways satisfied by the parent node i.e. the pathway set of a child node is a super-set of that of parent node.

3.2 Construction of Pathway Tree

Here, we describe the method of constructing the pathway tree. The conflict of a pathway with a node indicates conflict with its pathway set. The following steps are to be followed for construction of tree:

1. The initial contiguous set of non-conflicting pathways is added to the root node.

2. For each new pathway in the list, traverse the tree from root till it gets added to a node's list of pathways. Three cases arise here:
 - (a) If the node has two children with only one of them conflicting with the current pathway, set the non-conflicting node as current node. If both children are non-conflicting, choose either one with equal probability. Else, create a new node containing the parent's list of pathways and add the current pathway to its list and stop.
 - (b) If the node has only one child, and the child is conflicting, then proceed with creating a new node as mentioned earlier. Else, either choose the child as current node or create a new node (as mentioned earlier) with equal probability.
 - (c) If the node has no children, then create a new node (similar to previous step).

If a new node is created, the procedure terminates for the current pathway. Otherwise, steps (i), (ii) and (iii) are repeated for the new current node.

3.3 Selecting the optimum level of tree

Each level of the tree created, contains nodes representing a family of BNs and thus a PBN. The nodes in levels near the root contain fewer pathways while the ones at leaves have more pathways. Our goal is to strike a balance between them such that an optimum number of pathways, highly reliable according to information from time-series expression data are considered in our model. The time-series expression data gives us the regulatory genes (and their regulation: activation or inhibition) for each gene. At each level, we compute a score by summing the influences of these regulatory genes on BN across all target genes and across all nodes of that level. The score is then normalized by the number of nodes in that level. Let the number of nodes in i^{th} level be r_i and let n_{ik} denote the k^{th} node in i^{th} level. Let S_{kp}^i denote the set of regulatory genes for p^{th} gene in n_{ik} . $I_j \left(f_{n_{ik}}^{(p)} \right)$ denotes the influence of

gene x_j on the predictor function of p^{th} gene in n_{ik} . The score for level i is defined as,

$$\text{Score}(i) = \frac{\sum_{k=1}^{r_i} \sum_{p=1}^n \sum_{j|x_j \in S_{kp}^i} I_j(f_{n_{ik}}^{(p)})}{r_i} \quad (3.1)$$

The level with the highest score is selected as the optimum level, since it correlates the best with information from time-series expression data.

Thus, the PBN is constructed by using a linear combination of BN families of the optimum level and proportionally weighing each node by the size of its sub-tree.

Chapter 4

Performance Evaluation

In a n gene biological system, we randomly generate a set \mathbf{P} of non-conflicting pathways. Then we create m sets of pathways, each containing the set \mathbf{P} as subset, plus some additional pathways, non-conflicting with \mathbf{P} . Now, each of m sets of pathways represents a BN family. We mix them in a random proportion to generate our ground-truth PBN. Let p be the vector of coefficients of these m BNs in the PBN.

Selecting a random initial state and performing montecarlo simulations of the transition probability matrix of PBN gives us a time-series data of gene states (boolean values), which is then used to infer regulatory relationships. The pathway set used for constructing the tree contains the set \mathbf{P} followed by other pathways in those m sets. The algorithm constructs a tree, with h^{th} level being optimal. Let us assume the h^{th} level has r nodes. The normalized hamming distance metric for comparing two BN families is

$$\mu_{ham} = \frac{1}{n * 2^n} \sum_{i=1}^n \sum_{k=1}^{2^n} \left[f_i(x_k) \oplus f'_i(x_k) \right], \quad (4.1)$$

where $f_i(\cdot)$ and $f'_i(\cdot)$ represent Boolean functions of gene i in ground-truth and the inferred network ; x_k represents a binary state vector. The \oplus operator returns 0.5 in case either of the operand is a don't care, while the usual definition holds for other cases. However, in our case, the ground-truth network contains a weighted combination of m BN families. Let the reconstructed network contains r BN families. Accordingly, the distance metric for

comparison is as follows

$$\mu'_{ham} = \frac{\sum_{k=1}^r \left[\left(\sum_{i=1}^m \mu_{ham}(n_k, BN_i) * p(i) \right) * size(n_k) \right]}{\sum_{k=1}^r size(n_k)}, \quad (4.2)$$

where n_k denotes the k^{th} node in the optimum level h .

Chapter 5

Results and Discussion

The experiment described in the previous section is performed for $m = 5, 10, 15$ and $n = 4, 5, 6, 7$ genes. For each case, the measure is averaged over an ensemble of 100 biological systems, each containing a unique set of pathways and time-series expression data, shown in Table 5.2. The time-series contains 100 points and the value $l = 7$ was chosen for the inference algorithm. A sample output of the algorithm for $n = 6, m = 10$ including the tree, scores of different levels and the optimum level is shown in Fig. 5.1. Here, ' ABb ' denotes the pathway $A \xrightarrow{1:1,b} B$ and \mathbf{P} is the set of non-conflicting pathways. The corresponding regulatory relationships are shown in Table 5.1. The output PBN is obtained by switching between BN_1, \dots, BN_5 of optimum level 4 with probabilities $[\frac{2}{7}, \frac{2}{7}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}]$ respectively as shown in Fig 5.2.

The distance between ground-truth and the reconstructed PBN reduces with increase in number of genes. This is due to relatively more number of pathways in networks with lesser number of genes. Thus, less number of constraints leads to more efficient reconstruction. Increase in m also reduces the distance measure. This shows more amount of data results in better inference. However, the difference between $m = 10$ and $m = 15$ is less pronounced for $n = 6, 7$. This could be due to saturation in quality of new data available. In fact, presence of more inconsequential pathway information can only degrade the accuracy of the model.

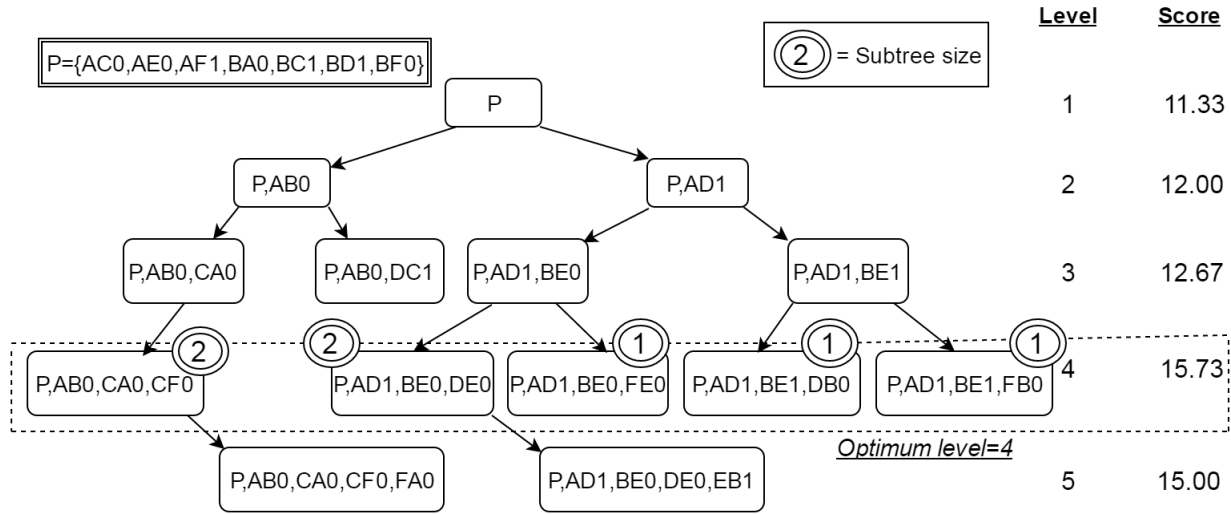


Figure 5.1: Pathway Tree for $n = 6, m = 10$

Table 5.1: Regulatory Relationships for $n = 6, m = 10$

a_{ij}	1	2	3	4	5	6
1	0	-1	0	0	0	0
2	0	0	0	-1	1	1
3	0	1	0	0	0	0
4	0	0	1	0	-1	0
5	1	0	0	-1	0	0
6	0	-1	1	0	-1	0

Table 5.2: Distance Measure between the ground-truth and reconstructed PBN

n	$m = 5$	$m = 10$	$m = 15$
	μ'_{ham}	μ'_{ham}	μ'_{ham}
4	0.524	0.458	0.428
5	0.444	0.371	0.357
6	0.344	0.281	0.278
7	0.275	0.255	0.243

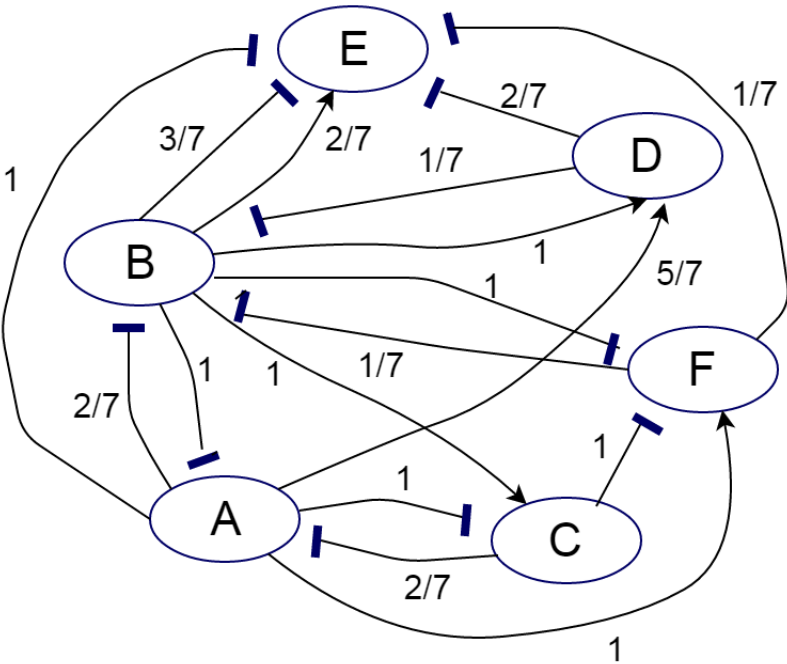


Figure 5.2: Output PBN for $n = 6, m = 10$

Chapter 6

Inference of Yeast Cell Cycle Network

The cell cycle is a vital biological process in which one cell grows and divides into two daughter cells. It consists of four phases, G1, S, G2, and M. Its regulation is highly conserved among eukaryotes [7]. From the 800 genes involved in cell cycle process of a budding yeast, Li et al. [5] constructed a network of 11 key regulators Cln3, MBF, SBF, Cln1, Clb5, Clb1, Mcm1, Cdc20, Swi5, Sic1, and Cdh1 which we shall refer to as A, B, C, D, E, F, G, H, I, J, and K respectively. We use the pathways and time-series expression data as given in [5] and shown in Table 6.1. We then compute the PBN using our algorithm. The resultant pathway tree is shown in Fig. 6.1. The output PBN switches between $BN_1, BN_2, BN_3, BN_4, BN_5,$ and BN_6 of optimum level 5 with probability $[\frac{1}{7}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}, \frac{2}{7}]$ respectively. The reconstructed PBN is shown in Fig. 6.2. For convenience, only pathways distinct with parent node are shown for each node. Our algorithm incorporates the essential pathways and gives proportional weightage to other pathways in modelling the cell cycle trajectory.

Table 6.1: Temporal evolution of state for yeast cell cycle

Time	Cln3	MBF	SBF	Cln1	Clb5	Clb1	Mcm1	Cdc20	Swi5	Sic1	Cdh1	Phase
1	1	0	0	0	0	0	0	0	0	1	1	Start
2	0	1	1	0	0	0	0	0	0	1	1	G1
3	0	1	1	1	0	0	0	0	0	1	1	G1
4	0	1	1	1	0	0	0	0	0	0	0	G1
5	0	1	1	1	1	0	0	0	0	0	0	S
6	0	1	1	1	1	1	1	0	0	0	0	G2
7	0	0	0	1	1	1	1	1	0	0	0	M
8	0	0	0	0	0	1	1	1	1	0	0	M
9	0	0	0	0	0	1	1	1	1	1	0	M
10	0	0	0	0	0	0	1	1	1	1	0	M
11	0	0	0	0	0	0	0	1	1	1	1	M
12	0	0	0	0	0	0	0	0	1	1	1	M
13	0	0	0	0	0	0	0	0	0	1	1	G1

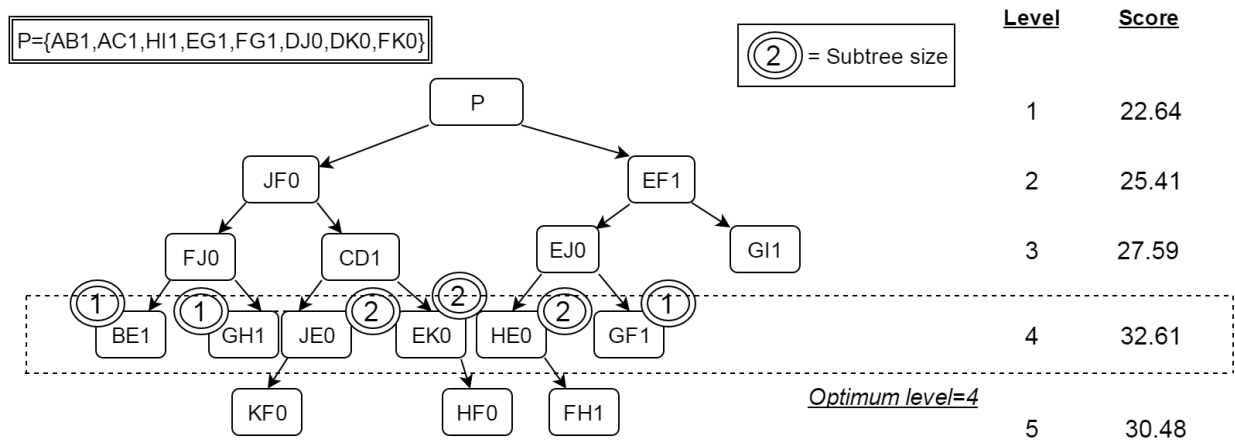


Figure 6.1: Pathway Tree for Yeast Cell Cycle

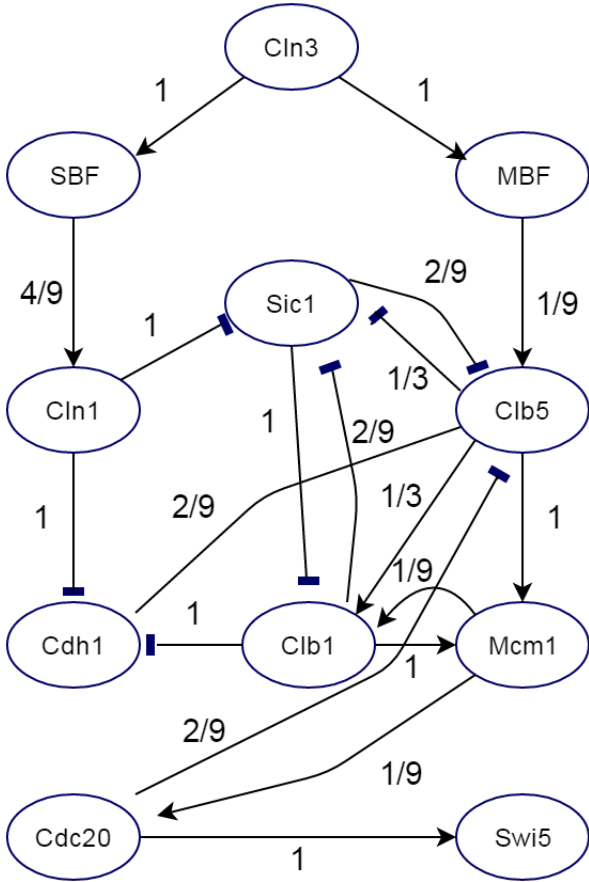


Figure 6.2: Output PBN for Yeast Cell Cycle

Chapter 7

Conclusion

In this project, we proposed an algorithm to infer a PBN for a biological system using biological pathways and time-series gene expression data. Pathways represent prior biological knowledge while time series data is obtained experimentally. The model space of PBN is huge compared to the amount of data available. Thus, a unique solution is impractical, given the fact that the data is noisy. Our solution overcomes this limitation by learning a model which makes systematic use of these two different forms of biological data. Future work will involve integration of multiple forms of such biological data to infer a more robust model.

Bibliography

- [1] Carlos HA Higa, Vitor HP Louzada, Tales P Andrade, and Ronaldo F Hashimoto. Constraint-based analysis of gene interactions using restricted boolean networks and time-series data. In *BMC proceedings*, volume 5, page 1. BioMed Central, 2011.
- [2] Sui Huang. Gene expression profiling, genetic networks, and cellular states: an integrating concept for tumorigenesis and drug discovery. *Journal of Molecular Medicine*, 77(6):469–480, 1999.
- [3] Stuart A Kauffman. Metabolic stability and epigenesis in randomly constructed genetic nets. *Journal of Theoretical Biology*, 22(3):437–467, 1969.
- [4] Ritwik K Layek, Aniruddha Datta, and Edward R Dougherty. From biological pathways to regulatory networks. *Molecular BioSystems*, 7(3):843–851, 2011.
- [5] Fangting Li, Tao Long, Ying Lu, Qi Ouyang, and Chao Tang. The yeast cell-cycle network is robustly designed. *Proceedings of the National Academy of Sciences of the United States of America*, 101(14):4781–4786, 2004.
- [6] Kevin Murphy, Saira Mian, et al. Modelling gene expression data using dynamic bayesian networks. Technical report, Technical report, Computer Science Division, University of California, Berkeley, CA, 1999.
- [7] Andrew Wood Murray and Tim Hunt. *The cell cycle: An Introduction*, volume 251. Oxford University Press New York, 1993.
- [8] Hongjia Ouyang, Jie Fang, Liangzhong Shen, Edward R Dougherty, and Wenbin Liu. Learning restricted boolean network model by time-series data. *EURASIP J. Bioinformatics and Systems Biology*, 2014:10, 2014.
- [9] Ilya Shmulevich, Edward R Dougherty, Seungchan Kim, and Wei Zhang. Probabilistic boolean networks: a rule-based uncertainty model for gene regulatory networks. *Bioinformatics*, 18(2):261–274, 2002.
- [10] Eugene P van Someren, Lodewyk FA Wessels, and Marcel JT Reinders. Linear modeling of genetic networks from experimental data. In *ISMB*, pages 355–366, 2000.