

Freie Universität zu Berlin



Masterthesis

Inference of Boolean Networks considering real-life time course Data

Nina Valery Kersten

Supervisors

Prof. Dr. Heike Siebert
Prof. Dr. Alexander Bockmayr

Advisor

Phd. Robert Schwieger

November 20, 2018

Abstract

The survival of a cell and eventually of its organism depends mostly on the reliable interaction between different kinds of substances. Different functionalities inside and outside a cell like proliferation, division and apoptosis are connected to different regulatory networks in a system. Small malfunction in these regulatory networks could cause a diseases from low impact for the organism to a big one. Learning these regulatory networks, its structure and dynamical behaviour is necessary to find a solution. Several effort have been made to reconstruct a regulatory network. In this work the focus is on Boolean networks which provide a simplified version of a system. In a Boolean network delivers binary information about the state of a substance, such that a gene is expressed or not or a protein's concentration is above or below a certain threshold. Here a real-life time course data set is used from a platform so called DreamChallenge to create a Boolean Network and compared to the best network created so far from a DreamChallenge Leaderboard group. The creation of the network is attended by preprocessing the data and showing which kind of influence the input data, the preprocessing and choice of the learning algorithm has to reconstruct a biological network reliably.

Contents

List of Figures	vii
List of Tables	ix
1 Introduction	1
2 Background	3
2.1 Biological Background	3
2.2 Preprocessing	8
2.3 Boolean Network,Interaction Graph and State Transition Graph	10
2.4 Network Evaluation	14
3 Materials and Methods	17
3.1 Inferencealgorithms	17
3.2 PyBoolNet	20
3.3 Data Selection	20
3.3.1 Example data set	20
3.3.2 HPN-DREAM breast cancer data set	20
4 Appilcation of the Pipeline	23
4.1 Example data set	23
4.2 HPN-DREAM breast cancer data set	23
5 Results	25
6 Conclusion	27
References	29

List of Figures

2.1	Signal Transduction	4
2.2	Transcriptional Gene Regulatory Network (GRN)	5
2.3	Direct binarization vs. Iterative binarization	9
2.4	Interaction Graph	12
2.5	Synchronous and Asynchronous State Transition Graph	14

List of Tables

3.1	20
3.2	20
3.3	20

1 Introduction

The development and functioning of a cell and an organism in general is a product of a complex cellular machinery. This machinery is compound by the interaction of genes, proteins, mRNA and many other substances to induce a cascade of extracellular signals transducted by mechanisms of the cell membrane, reaching the nucleus of the cell, initiating a transcription process that controls the production and abundance of proteins. Proper functioning of these networks is essential to the survival and adaption of all living organisms, while malfunctioning of these networks has been identified as the cause of various diseases [?]. To understand the behaviour of a biological system it is necessary to find and analyze the main important processes in a system in a dynamical manner. Therefore high-throughput techniques provide a big abundance of information about various biological interactions measured over a series of time. Biological information can be considered as different systems such as signal transduction, gene regulation, protein-protein-interaction or metabolism. These information are put in a network which can be yielded by several strategies like Bayesian networks, Boolean regulatory networks, Ordinary differential equation models and Neural networks[Saadatpour & Albert, 2013]. Once a network is constructed further analysis of the network by validating the network (e.g. perturbations like gene manipulation and external treatments, reductions etc.) can be done to figure out the main interactions whose disfunctionalities cause diseases or structural and steady-state property analysis. This masterthesis focuses on constructing a pipeline by creating a boolean network from a real-life time course data set. In the following section the biological background of data, it's preprocessing by normalization and discretization, the graphtheoretical background and scoring methods are described. In the second section it is getting more detailed by giving an overview of different inference algorithms, describing a tool called PyBoolNet and showing with an example data set and a real-life time course data set of breast cancer celllines how the pipeline is applied. In the last section the scoring results for the real-life time course data set are analyzed reflecting the performance of the constructed network.

2 Background

2.1 Biological Background

Depending on the aim of a network inference the biological input data can be depicted by the interaction of proteins, genes and metabolic substances. In this section the intention of using different types of input data is explained and what potentially will be the occurring problems. Different types of biological input data provide different structures of the input data for further network inference algorithms. The main interactions of interest are the Gene Regulatory Networks (GRN) and the Protein-Protein Interactions (PPI) described below. With this interactions the most most complex cellular process of metabolic interactions can be analysed. Connections between biochemical reaction via substrate and product metabolites create complex metabolic networks. The focus is set on the different aspects of enzyme chemistry, enzyme structure and metabolite structure. Thus an individual's metabolism is determined by one's genetics, environment, and nutrition. By investigating the chemical structure of metabolites and systematically classify the functions of the enzymes the understanding of a metabolism and the prediction of enzyme function and novel metabolic pathways is improved.[?][?] Furthermore GRN and PPI can help to analyze the aspects of signal transduction. In the developement of multicellular organisms the action of extracellular growth factors activate a cascade of intracellular signaling pathways. These pathways regulate major aspects of cell regulation like cell proliferation, cell migration, cell differentiation, cell survival and cell death. To understand the developement of diseases (e.g. cancer) major processes (e.g. phosphorylation, ubiquitylation, methylation, etc.) of signal transduction pathways can be delighted by the prediction of a network. The concentration of signalling pathway underlies high fluctuation over time due to transcriptional and translational regulation, such that the inference of a network is a challenging task [Kestler et al., 2008]

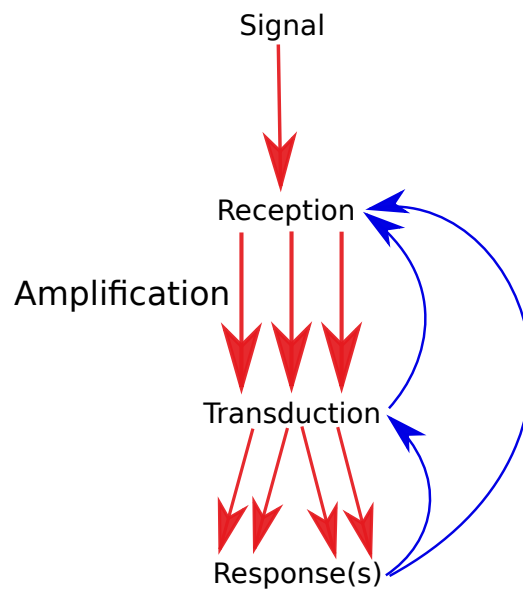


Figure 2.1: Signal Transduction. An environmental signal (e.g. hormone) interacts with a cellular component, most often a cell-surface receptor. The information that the signal has arrived is then converted into other chemical forms, or transduced. The signal is often amplified before introducing a response. Feedback pathways regulate the entire signaling process.[????]

Transcriptional Gene Regulatory Networks

In a Transcriptional Gene regulatory network (GRN) the nodes are depicted by the genes and the arcs are directed and show whether a gene produces RNA (transcript of the source gene, resp. regulator) which inhibits oder activtes the target gene (regulatee). For network inference computational algorithms take the mRNA expression levels of genes as the input data.

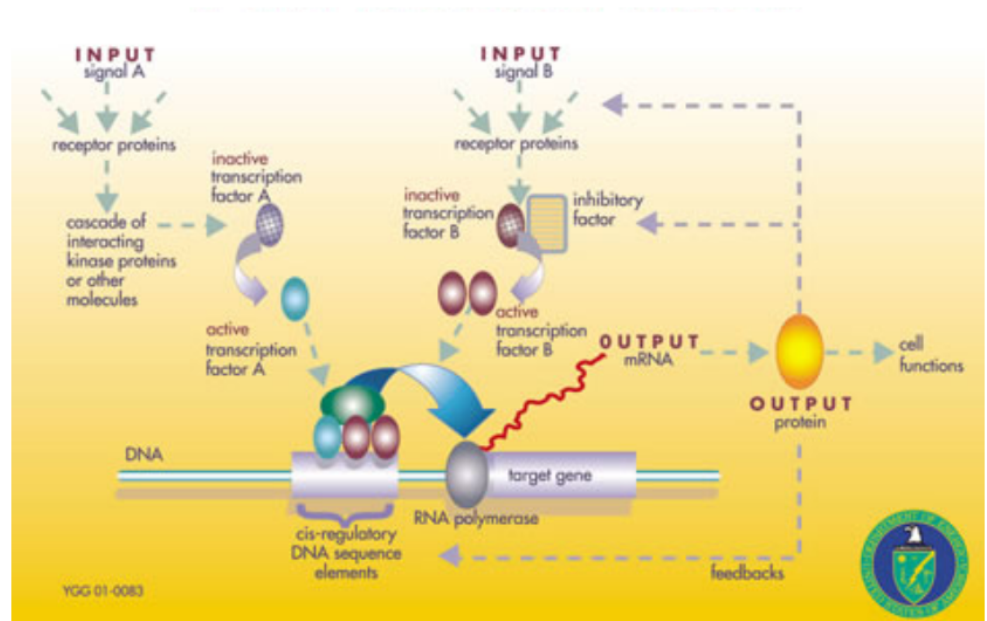


Figure 2.2: Transcriptional Gene Regulatory Network (GRN).

In this example two different signals have an impact of a single target gene. Signal molecule A triggers the conversion of inactive transcription factor A (green oval) into an active form that binds directly to the target gene's cis-regulatory sequence. The process for signal B is more complex. Signal B triggers the separation of inactive B (red oval) from an inhibitory factor (yellow rectangle). B is then free to form an active complex that binds to the active A transcription factor on the cis-regulatory sequence. The net output is expression of the target gene, leveled by A and B. Thus cis-regulatory DNA sequences with the proteins that assemble on them, integrate information from multiple signaling inputs to produce mRNA-Output . [?]

Protein-Protein-Interaction

In contrast to the gene regulatory interaction network the protein-protein interactions (PPIs) act directly among themselves. Thus the nodes in a network are the interacting proteins. Proteins interact by physical contacts (e.g. electrostatic forces) of high specificity. PPIs play a big role in electron transfer, signal transduction, transport across membranes and cell metabolism. A variety of techniques are known to detect PPIs. The most applied ones are immuno-precipitations and the yeast two-hybrid approach. The two-hybrid assay is not a reliable indication that two proteins interact *in vivo*, because the two interacting proteins are overexpressed. Thus the interaction may not be present in the wild type cells where the concentrations may be significantly lower. For this reason additional information are included to figure out the occurrence of true interaction, such as cellular localization and mRNA expression level. The underlying assumption is that true interactions are likely to occur between proteins involved in the same biological process, proteins found in the same cell compartment, and proteins whose mRNA are coexpressed. Interaction networks of PPIs may depict how drug-protein interactions lead to toxic side effects. [Pellegrini et al., 2004]

The collection of PPI data is done by a technique so called Reverse phase protein lysate microarray (RPMA). This technique is divided up into 6 parts. First starting with the sample collection. An inhibitor or stimulus in form of drugs is added to a set of cell lines at the same time and the cell lines are then processed at different time points. Secondly in the Cell Lyses step cell fragments are lysed with a cell lysis buffer to obtain high protein concentration. The choice of a buffer decides the quantity of proteins can be lysed out of the cell. Afterwards cell lysed probes are diluted. In the Antibody screening the lysates are pooled and resolved by SDS-PAGE followed by western blotting on a nitrocellulose membrane. The membrane is cut into 4mm strips. Each slide is probed with a different antibody, primary with a secondary antibody. For fluorometric detection a primary and secondary antibody are diluted. Detection reagent is put on each slide. Signal amplification and detection is done by an optical flatbed scanner if colorimetric technique is used or by laser scanning. Subsequently the data set structure is determined by deleting missing data points and outliers from the set. Then the data set is normalized.

The strength of RPMA is the high throughput, ultra-sensitive detection of proteins from extremely small numbers of input material which is not possible for western blotting and ELISA. The small spot size on the microarray, ranging in diameter from 85 to 200 micrometres, enables the analysis of thousands of samples with the same antibody in one experiment. The

high sensitivity of RPMA allows for the detection of low abundance proteins or biomarkers such as phosphorylated signaling proteins from very small amounts of starting material such as biopsy samples, which are often contaminated with normal tissue. A great improvement of RPMA over traditional forward phase protein arrays is a reduction in the number of antibodies needed to detect a protein. The protein isn't detected directly which helps to preserve the proteins. Antibodies, especially phospho-specific reagents, often detect linear peptide sequences that may be masked due to the three-dimensional conformation of the protein. This problem is overcome with RPMA as the samples can be denatured, revealing any concealed epitopes.

The weakness of RPMA are batch effects caused by the choice of the right buffer, quantity of the proteins and the antibody performance. The choice of the right buffer decides the quantity of proteins which can be lysed out of the cell. Little or poor quality of starting material and a long storage time causes low protein. It might be useful to improve the antibody performance by validating it with a smaller sample size under identical conditions before starting with the actual sample collection. Currently the number of signaling proteins for which antibodies exist to get an analyzable signal is small.

All these facts should be considered in later analysing steps.

2.2 Preprocessing

After generating the biological data some preprocessing like normalisation and discretization has to be done. Therefore several strategies are known. Normalization is an essential step, because data can contain outlier, the abundance of some proteins or mRNA is often higher than others and obtaining biological data from the lab could cause several batch effects. In 3.3.2 HPN-DREAM breast cancer data set a normalization method is described more detailed. Before starting inferring a boolean network from real-life time-series data the data has to be discretized such that each value (e.g. concentration measurement) measured at a certain time point of a particular substance (e.g. gene, protein) has either a value 1 or 0. The choice of the appropriate discretization algorithm decides about accuracy of a network. In this part of the chapter three discretization algorithms are introduced.

Two clusters k-means binarization

The time-series data is divided into two clusters directly. One cluster contains all the values with the higher mean being set to 1. In the other cluster all the values with the lower mean being set to 0 are combined. This binarization strategie is fast and simple but may exclud some essential information like about oscillations and fluctuations.[Berestovsky & Nakhleh, 2013]

Iterative k-means binarization

A depth d of clustering is set followed by a set of initial number of cluster $k = 2^d$. The input consists of tim-series data $S = \{S_1, \dots, S_n\}$ of n species (e.g.different genes or proteins),each of size $m + 1$, where $S_i(t) \in \mathbb{R}^+$ ($0 \leq t \leq m$) is the concentration of species i at time t . In each iteration the data is classified into k dijoint clusters $C_{s_i}^1, \dots, C_{s_i}^x$. In each Cluster all its values are replacd by the clusters mean $\mu(C_{s_i}^x)$. Here $d = 3$ in the beginning until $d = 1$ such that the two cluster k-means binarization method can be applied.

Example 2.1. Assume we have a time-series data with measurements for a gene A . Starting with a depth of $d = 3$ we have initally $k = 8$ clusters for each gene. Resulting in $\{C_{s_A}^1, \dots, C_{s_A}^8\}$, each cluster containing values of time-series data for A . Now for each cluster the mean $\{\mu(C_{s_A}^1), \dots, \mu(C_{s_A}^8)\}$ is calculated. Afterwards values in a cluster are replaced by its mean. This is done 2 times more such that we get the final cluster with $d = 1$.

In Figure 2.3 the advantages of an iterative binarization in contrast to the direct binarization is shown. One of the big advantages is that in iterative clustering the information about oscillations is maintained.

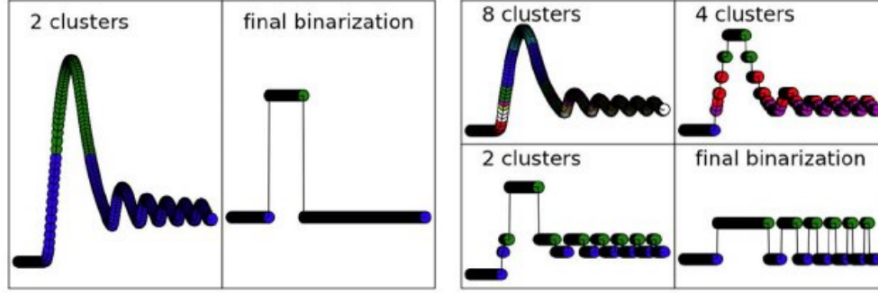


Figure 2.3: Direct binarization vs. Iterative binarization. More detailed binarization is achieved with higher values of d

BASC A -Binerization

The BASC A binerization algorithm is a bit more complex, thus here a short description is given. For more details have a look in the paper The algorithm is divided up into three parts, starting with computing a series of stepfunctions, then finding the strongest continuity in each step function and estimate the location and variation of the strongest discontinuities. Firstly an intial step function is obtained by rearranging the original time-series measurements in increasing order. Step functions with fewer discontinuities are calculated such that each minimizes the eucledian distance to the initial step function. Afterwards the strongest discontinuity in each step function is found by a high jump size (derivative) in combination with a low approximation error. Finally time-series measurements of gene expression values can be excluded from the network based on the estimations of the location and variation of the strongest discontinuities.

2.3 Boolean Network, Interaction Graph and State Transition Graph

In this section the knowledge of the discretized biological data is put into a graphtheoretical context of a boolean network. Here mathematical definitions are given, explained by an example.

Definition 2.2. Undirected Graph and Directed Graph

*In general, an undirected graph $G = (V, E)$ is defined as a set of vertices V describing the nodes of the system and a set of undirected edges $E = \{(i, j) | i, j \in V\}$ that define a relationship between node i and j . While in a **directed graph** $G = (V, A)$ is an ordered pair, defined as a set of vertices V (nodes) and a set of directed edges A (arcs). A set of directed edges $A = \{(i, j) | i, j \in V\}$ describes the flow of information in network, where (i, j) describes the flow from i (tail) to j (head).*

Definition 2.3. Boolean Network

A boolean network is a directed graph $G(X, F)$, defined by a set of nodes X in a binary vector $X(t) = \{x_1(t), \dots, x_n(t)\}$ representing state of a system at time t . Each element $x_i \in X$ corresponds to the state $x_i = 1$ or $x_i = 0$ of a species i . A set $F = \{f_1, \dots, f_n\}$ of n transition functions contains a particular function for each x_i .

For every $f_i \in F$ s.t. $1 \leq i \leq n$,

$$f_i(X(t)) = x_i(t + 1)$$

, where $f_i(X(t))$ defines the next state of x_i at time t in the network.

Definition 2.4. 2.Definition:BooleanNetwork *A network (X, F) is Boolean iif all its components are Boolean. In Boolean networks, every update function $f \in F$ is therefore a n -variable Boolean function $f : \mathbb{B}^n \leftarrow \mathbb{B}$ which we can represent by a Boolean expression over the n input variables V .*

[?]

This means a given boolean network consists of n nodes (resp. variables, species, vertices) representing the components of a system and the directed edges (resp. arcs) representing the orientation of interaction between the nodes. Thus a regulator (resp. source node) has an impact on the regulatee (resp. target node). Each node can have two possible initial states $x_i \in \{0, 1\}$ describing the its activity. The activity of a node means a qualitative rate a gene is being transcribe $x_i = 1$ or not $x_i = (0)$, a transcription factor is active oder

inactive, a protein's concentration is above or below a certain threshold (e.g. phosphorylated or un-phosphorylated). The future state of a node is determined by boolean rule so called a transition function $f_i(X(t))$. [Berestovsky & Nakhleh, 2013] [?]

An interaction graph is some kind of an abstraction of the boolean network represented in a directed graph $IG(X, A)$ where the directed edges have a source node x_i and a target node x_j . The edges can be weighted positive or negative by an edge weight assigned to each directed edge by $\omega : A \rightarrow \{+, -\}$.

Definition 2.5. Interaction Graph

The interaction graph of (X, F) is the directed graph (X, \leftarrow) that consists of the node set X and the ars set $\leftarrow = \leftarrow_F \subseteq X \times X$ with $(y, x) \in \leftarrow$ iff f_x depends on y .

[?]

In a biological system genes (resp. proteins) may have a positive or negative influence on other genes (resp. proteins). This is represented by the labeled edges in an interaction graph. Logical operator AND, OR, NOT (resp. $\&\&$, $\|$, $!$; resp. \wedge , \vee , \neg) are commonly used to describe these interactions. For instance a gene x_A could be influenced by another gene x_B positive and by a second gene x_C negative. Then the boolean algebra would look like this:

$$x_A = x_B \text{ AND NOT } x_C$$

$$x_A = x_B \ \&\& \ !x_C$$

$$x_A = x_B \wedge \neg x_C$$

Furthermore the dynamics of a system can be simulated by repeatedly applying the transition functions and updating the "current" state in a synchronous or asynchronously manner. This yields from an interaction graph to a state transition graph. [Saadatpour & Albert, 2013]

A transition function specifies how F determines for every $x \in S$ a set of states that are reachable by a single transition. The resulting relation between the states of a model is usually thought of as a directed graph (S, \rightarrow) called the state transition graph (STG). The STG is the basis for analysis and for simulations of the dynamics of a qualitative model. [?]

Definition 2.6. State Transition Graph

A state transition graph (STG) is a directed graph with a set of nodes represented by a set of binary vectors $F(t+1) = \{f_1(X(t+1)), \dots, f_n(X(t+1))\}$ representing the updated states of all variables after all of the functions in F have executed. The arcs denote possible transitions from one binary state vector to another.

[?] [Lee et al., 2002].

In a synchronous state transition graph the states of all variables are updated simultaneously after all of the functions in F have executed. In an asynchronously updated STG the states are updated one at a time by randomly choosing a transition function $f_i \in F$ and updating the state of x_i immediately.

Example 2.7. In this Example the processing of a boolean network is shown, starting with the construction of an interaction graph and later converted into a state transition graph with a synchronous and an asynchronous update.

The first figure shows an interaction graph of a boolean network with $X = \{x_1, x_2, x_3\}$ three nodes and four activating arcs (black arrow) and two inactivating arcs (red arrow).

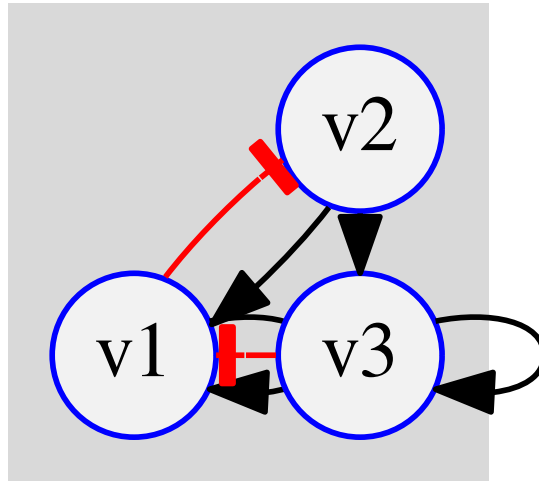


Figure 2.4: Interaction Graph.

From the interaction graph the set of transition functions F can be created such that the state transition graph can be calculated.

$$f(x_{1,2,3}) = \begin{pmatrix} x_1 & \wedge x_2 & \wedge \neg x_3 \\ \neg x_1 & & \\ & x_2 & \wedge x_3 \end{pmatrix}$$

For every possible state of $x_i \in X$ the next state $f_i(X(t)) = x_i(t+1)$ is calculated shown in the computation (2.1)-(2.8) below. This computation provide the new states for a synchronously updated state interaction graph displayed in the left graph in Figure 2.5. [Allisabeth Remy et al., 2008]

$$x(t) = (0, 0, 0) \rightarrow f(x(t+1)) = (0, 1, 0) \quad (2.1)$$

$$\textcolor{red}{x(t) = (0, 0, 1)} \rightarrow \textcolor{red}{f(x(t+1)) = (0, 1, 0)} \quad (2.2)$$

$$x(t) = (0, 1, 0) \rightarrow f(x(t+1)) = (0, 1, 0) \quad (2.3)$$

$$x(t) = (1, 0, 0) \rightarrow f(x(t+1)) = (0, 0, 0) \quad (2.4)$$

$$x(t) = (1, 1, 0) \rightarrow f(x(t+1)) = (1, 0, 0) \quad (2.5)$$

$$\textcolor{red}{x(t) = (1, 0, 1)} \rightarrow \textcolor{red}{f(x(t+1)) = (0, 0, 0)} \quad (2.6)$$

$$x(t) = (0, 1, 1) \rightarrow f(x(t+1)) = (0, 1, 1) \quad (2.7)$$

$$\textcolor{red}{x(t) = (1, 1, 1)} \rightarrow \textcolor{red}{f(x(t+1)) = (0, 0, 1)} \quad (2.8)$$

The asynchronous updated state transition graph is computed by adding intermediate computation steps to the synchronus computation. Therefore the red highlighted computation are split by the amount of state changes. For instance in computation (2.2) $x(t) = (0, 0, 1)$ has two updates of states, x_2 changes and x_3 changes, too. As we know from the descriptin of asynchronous STG's in biological systems processes happen uncommonly at the same time. Thus the initial state $x(t) = (0, 0, 1)$ provides two updates: $f(x(t+1)) = (0, 1, 1)$ and $f(x(t+1)) = (0, 0, 0)$. The same procedure is done with (2.6) and (2.8) and finally the asynchronous STG can be drawn, shown by the right graph in Figure 2.5.

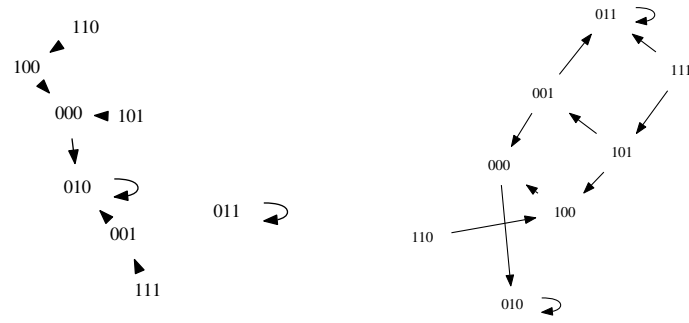


Figure 2.5: Synchronous and Asynchronous State Transition Graph. Left: Synchronous State Transition Graph; Right: Asynchronous State Transition Graph

2.4 Network Evaluation

Evaluating a boolean network is done by the Area under the Curve (AUC) of the Receiver-Operating-Characteristic-Curve (ROC) by comparing the interaction graph of a training data set against an interaction graph of a goldstandard data set. This score returns a value between $[0, 1]$. It is desirable to get a value close to 1 which means the model predicted the network quite well. Calculating the ROC provides a set of additional information like the accuracy, precision and recall.

For every scoring type a specific case is used described in the table below.

True Positive (TP)	False Positive (FP)
False Negative (FN)	True Negative (TN)

Definition 2.8. Accuracy. Fraction of predictions our model got right.

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{TP + TN}{TP + TN + FP + FN} \quad (2.9)$$

Definition 2.9. Precision. Returns the proportion of positive identifiers which was accurately correct.

$$Precision = \frac{TP}{TP + FP} \quad (2.10)$$

Definition 2.10. Recall. Returns the proportion of actual positives identified correctly.

$$\boxed{Recall = \frac{TP}{TP + FN}} \quad (2.11)$$

For determine the Receiver-Operating-Characteristic-Curve a True-Positive-Rate (TPR) and a False-Positive-Rate (FPR) is calculated.

Definition 2.11. True-Positive-Rate (TPR). The TPR values are for the y-axis of the ROC.

$$\boxed{TPR = \frac{TP}{TP + FN}} \quad (2.12)$$

Definition 2.12. False-Positive-Rate (FPR). The FPR values are for the x-axis of the ROC.

$$\boxed{FPR = \frac{FP}{FP + TN}} \quad (2.13)$$

Accuracy is not always an appropriate scoring method, because, assigning every object to a larger set achieves a high proportion of correct predictions, but is not generally a useful classification. For this reason Blanced Accuracy and the Mathew Correlation Coefficient are introduced.

Definition 2.13. Balanced Accuracy (BACC). Fraction of predictions our model got right divided by 2.

$$\boxed{BACC = \frac{\frac{\text{Number of correct predictions}}{\text{Total number of predictions}}}{2} = \frac{\frac{TP+TN}{TP+TN+FP+FN}}{2}} \quad (2.14)$$

Definition 2.14. Matthew Correlation Coefficient (MCC). The MCC measures the quality of binary classifications and is a correlation coefficient between observed and predicted binary classifications which returns a value between -1 and 1 ; $MCC \in [-1, 1]$. Where a value close to 1 denote a perfect prediction, a value close to 0 means that the prediction is not better than a random one and a value close to -1 describes a total disagreement between prediction and observation.

$$\boxed{MCC = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}} \quad (2.15)$$

Definition 2.15.

Example 2.16. Given is model that classified 100 tumors as malignant (the positive calss) or benign (the negative class):

True Positive (TP):	False Positive (FP):
False Negative (FN):	True Negative (TN):

3 Materials and Methods

3.1 Inference algorithms

BIBN (Bayesian inference approach for a Boolean network) REVEAL (Reverse Engineering algorithm) PCA-CMI (Path consistency algorithm- Conditional mutual information) ARCANE (Time-delay algorithm for reconstruction of accurate cellular networks) MIDER (Mutual information distance and entropy reduction)[?] defines a mutual information based distance between genes to specify the directionality BESTFIT ()

These mutual information-based methods are computationally expensive, because they are implemented to compute exact mutual information values over all possible combinations of genes.

RelNet (Relevance network algorithm)

CLR (context likelihood of relatedness method) CST (chi-square test) [Chai et al., 2014]

Boolean Model

Boolean models do not need any information about kinetic parameters, provide an important insight into the dynamics and the structure of a network. The relationships in a boolean network can be derived from a relatively small dataset. Furthermore a boolean model could make qualitative predictions of large complex networks more feasible Boolean Models are either probabilistic (PBN) or deterministic (DBN). In a PBN the next state of a gene is determined by a transition function f_i selected with a certain probability from a set of transition function F . Here the focus is on a deterministic boolean model.

Redundancy Removal

Real-life time course data could contain redundant information about steady-states and the significance of a transition. A steady-state is a point attractor which is a pair of equal consecutive states indicating a true-steady state iif. it was the last pair in the series. Data measured in a very fine time-scale is often giving false indication of steady-state signal. To overcome this false-steady-state-transition-problem each maximal consecutive sequence of identical states are all removed except for one of the states.

Regarding the significance of a transition the average number of bits needed is determined. If the average number of bits is above 1, this reduction skips informative transitions, which could be important in the learning step.

REVEAL

REVEAL (REVerse Engineering ALgorithm) is an inference algorithm which uses a deterministic transition table to infer boolean relationships between variables. After maximal 2^n iterations of the algorithm a "steady-state" (resp. point attractor) should be found with represent the boolean rule (transition function). REVEAL deals with calculation of individual's entropy in combination with the joint entropy and the mutual information.

Definition 3.1. Shannon-Entropy

The Shannon-Entropy is the probability of observig a particular symbol of event $p(x)$, within a given sequence.

$$H = - \sum p(x) \log p(x) \quad (3.1)$$

Example 3.2. Here $p(x)$ is the probability of observing value $x \in \{0, 1\}$ for a variable x_i .

x	0	1	1	1	1	1	1	0	0	0
y	0	0	0	1	1	0	0	1	1	1

$$H(x) = -p(0) * \log[p(0)] - [1 - p(0)] * \log[1 - p(0)] \quad (3.2)$$

$$H(x) = -0.4\log(0.4) - 0.6\log(0.6) = 0.97(40\%0, 60\%1) \quad (3.3)$$

$$H(x) = -0.5\log(0.4) - 0.5\log(0.6) = 1.00(50\%0, 50\%1) \quad (3.4)$$

An element x can have two possible states 1 (on) or 0 (off). H reaches its maximum when both possible states equally probable $H_{max} = \log(2) = 1$. Beside the individual entropy of x and y now the combined entropy is consulted.

Definition 3.3. Joint Entropy The joint entropy is defined by the pobability of occurences that x and y occur dependend on each other.

$$H(x, y) = - \sum p(x, y) \log p(x, y) \quad (3.5)$$

Example 3.4. The co-occurences of 1 and 0 in x and y are displayed in a quadratic matrix. Afterwards the joint entropy $H(x, y)$ is calculated in (3.6).

$$\begin{array}{c}
 \mathbf{y} \\
 \begin{array}{cc}
 1 & \boxed{\begin{array}{|c|c|} \hline 3 & 2 \\ \hline \end{array}} \\
 0 & \boxed{\begin{array}{|c|c|} \hline 1 & 4 \\ \hline \end{array}} \\
 0 & 1 \\
 \mathbf{x}
 \end{array}
 \end{array}$$

$$H(x, y) = -0.1\log(0.1) - 0.4\log(0.4) - 0.3\log(0.3) - 0.2\log(0.2) = 1.85 \quad (3.6)$$

In the last computational step the mutual information is calculated by the combination of the Shannon-Entropy and the Joint-Entropy.

Definition 3.5. Mutual Information The mutual information describes the rate of transmission.

$$\boxed{M(X, Y) = H(X) + H(Y) - H(X, Y)} \quad (3.7)$$

This equation can be extended n-times, for n nodes in a network.

$$\boxed{M(X, [Y, Z]) = H(X) + H(Y, Z) - H(X, Y, Z)} \quad (3.8)$$

The smallest subset x' that yields $M(x_i, x'_i)/H(x_i) = 1$ reflect the set of nodes (resp. genes, proteins) whose states determine the next state of the gene represented by a variable x_i .

Example 3.6. With the knowledge about Shannon entropy, joint entropy and the mutual information the boolean rules (resp. transition functions) for the example of state transition tables with the node set $n = A, B, C$ can be calculated.

Transition table "B":

input	at	time (t)	input	at	time ($t + 1$)
A	B	C	A'	B'	C'
0	0	0	0	0	0
0	0	1	0	1	0
0	1	0	1	0	0
0	1	1	1	1	1
1	0	0	0	1	0
1	0	1	0	1	1
1	1	0	1	1	1
1	1	1	1	1	1

Input entropies

H(A)	1.00
H(B)	1.00
H(C)	1.00
H(A,B)	2.00
H(B,C)	2.00
H(A,C)	2.00
H(A,B,C)	3.00

→

Determine the mutual information for A		
H(A') 1.00		
H(A',A) 2.00	M(A',A) 0.00	M(A',A)/H(A') 0.00
H(A',B) 1.00	M(A',B) 1.00	M(A',B)/H(A') 1.00
H(A',C) 2.00	M(A',C) 0.00	M(A',C)/H(A') 0.00

Table 3.2

Table 3.1

If $M(A', X) = H(A')$ then $M(A', X)/H(A') = 1$, then X exactly determines A' . This is here the case for B in A' , where A' denotes the output's state shown in the red highlighted line in Table 3.2. The iteration of REVEAL stops here and the boolean rule (resp. transition function) can be inferred.

input	output
B	A
0	0
1	1

→

$$f_A = B$$

Table 3.3

For more details the is referred to the paper [TS2B [REVEAL-PAPER]]

BESTFIT

FULLFIT

Error Assesement

3.2 PyBoolNet

3.3 Data Selection

3.3.1 Example data set

3.3.2 HPN-DREAM breast cancer data set

Now it is shown how the Pipeline can be applied to a real-life time course data set.

What is the Dream Challenge?

For a Boolean network inference the data of a platform so-called Dialogue on Reverse Engineering Assessment and Methods (DREAM) - Challenge is used. The DREAM-Challenge is a non-profit, collaborative community effort consisting of contributors from across the research spectrum of questions in biology and medicine. This organization was built in 2006 and publishes crowdsourcing challenges with transparent sharing of data, thus everyone can participate the challenge. The DREAM-Challenge has partnered with Sage Bionetworks, which provide the infrastructure by Sage Bionetworks Synapse platform to get access to the open collaborative data analysis. Overall the DREAM-Challenge is a helpful instrument to get real-life data, comparing results and interact with other researchers all over the world, while contribute solutions to biological and medical questions.[?]. The challenging question was to decide which Dream Challenge data set could be useful for this masterthesis. For inferring a Boolean network and further analysis of the state transition graph the desired data set should contain measurements of experiments with less perturbational information and a in a time course context. The Dream Challenge 5 dealing with gene-gene interaction,providing test and training data sets of gene expressions seemed to be an appropriate candidate. But there was less time course information and a high abundance of perturbation. Thus the Dream8 Challenge is was chosen. This challenge describes protein-protein interaction and measurements for multiple timepoints.

DREAM8

Data Collection

The collection of the HPN8 breast cancer PPI data is done by a technique so called Reverse phase protein array(RPPA). This technique is divided up into 6 parts:

Sample collection An inhibitor or stimulus in form of drugs is added to a set of celllines at the same time and the celllines are then processed at different time points.

Cell Lysis Cell fragments are lysed with a celllysis buffer to obtain high protein concentration.The choice of a buffer decides the quantity of proteins can be lysed out of the cell.

Dilution Dilution of the celllysed probes.

Antibody screening The lysates are pooled and resolved by SDS-PAGE followed by western blotting on a nitrocellulose membrane. The membrane is cut into 4mm strips. Each slide is probed with a different antibody, primary with a secondary antibody.

Fluorometric detection Primary and secondary antibody are diluted. Detection reagent is put on each slide. Signal amplification and detection is done by an optical flatbed scanner if colorimetric technique is used or by laser scanning.

Data set structure Missing data points and outliers are detected and deleted from the data set. The data set is normalized

4 Appilcation of the Pipeline

4.1 Example data set

4.2 HPN-DREAM breast cancer data set

5 Results

5.1

5.2

5.3

6 Conclusion

References

- Berestovsky, N., & Nakhleh, L. (2013, 06). An evaluation of methods for inferring boolean networks from time-series data. *PLOS ONE*, 8(6), 1-9. Retrieved from <https://doi.org/10.1371/journal.pone.0066031> doi: 10.1371/journal.pone.0066031
- Chai, L. E., Loh, S. K., Low, S. T., Mohamad, M. S., Deris, S., & Zakaria, Z. (2014). A review on the computational approaches for gene regulatory network construction. *Computers in Biology and Medicine*, 48, 55 - 65. Retrieved from <http://www.sciencedirect.com/science/article/pii/S0010482514000420> doi: <https://doi.org/10.1016/j.compbiomed.2014.02.011>
- De Las Rivas, J., & Fontanillo, C. (2010, 06). Protein-protein interactions essentials: Key concepts to building and analyzing interactome networks. *PLOS Computational Biology*, 6(6), 1-8. Retrieved from <https://doi.org/10.1371/journal.pcbi.1000807> doi: 10.1371/journal.pcbi.1000807
- Elena S. Dimitrova, J., M. Paola Vera Licona. (2010). Discretization of time series data. *Journal of Computational Biology*(1), 853 - 868. Retrieved from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3203514/> doi: <http://doi.org/10.1089/cmb.2008.0023>
- Kestler, H. A., Wawra, C., Kracher, B., & Kuehl, M. (2008). Network modeling of signal transduction: establishing the global view. *BioEssays*, 30(11-12), 1110–1125. Retrieved from <http://dx.doi.org/10.1002/bies.20834> doi: 10.1002/bies.20834
- Lee, T. I., Rinaldi, N. J., Robert, F., Odom, D. T., Bar-Joseph, Z., Gerber, G. K., ... Young, R. A. (2002). Transcriptional regulatory networks in *saccharomyces cerevisiae*. *Science*, 298(5594), 799–804. Retrieved from <http://science.sciencemag.org/content/298/5594/799> doi: 10.1126/science.1075090
- Pellegrini, M., Haynor, D., & Johnson, J. M. (2004). Protein interaction networks. *Expert Review of Proteomics*, 1(2), 239-249. Retrieved from <https://doi.org/10.1586/14789450.1.2.239> (PMID: 15966818) doi: 10.1586/14789450.1.2.239

- Saadatpour, A., & Albert, R. (2013). Boolean modeling of biological regulatory networks: A methodology tutorial. *Methods*, 62(1), 3 - 12. Retrieved from <http://www.sciencedirect.com/science/article/pii/S1046202312002770> (Modeling Gene Expression) doi: <https://doi.org/10.1016/j.ymeth.2012.10.012>
- Elisabeth Remy, Ruet, P., & Thieffry, D. (2008). Graphic requirements for multistability and attractive cycles in a boolean dynamical framework. *Advances in Applied Mathematics*, 41(3), 335 - 350. Retrieved from <http://www.sciencedirect.com/science/article/pii/S0196885808000146> doi: <https://doi.org/10.1016/j.aam.2007.11.003>

Hiermit erkläre ich an Eides statt, dass ich die vorliegende Arbeit selbständig und ohne unerlaubte fremde Hilfe angefertigt, andere als die angegebenen Quellen und Hilfsmittel nicht benutzt und die den benutzten Quellen und Hilfsmitteln wörtlich oder inhaltlich entnommenen Stellen als solche kenntlich gemacht habe.

Berlin, den 07. September 2018

.....
(*Unterschrift*)