

Linköping Electronic Articles in
Computer and Information Science
Vol. 3(2001): nr 7

Approaches for Simulating and Modeling Cell Regulation: Search for a Unified View Using Constraints

Jacques Cohen
Department of Computer Science
and
Center of Complex Systems
Brandeis University
Waltham, MA 02454
email: jc@cs.brandeis.edu

Linköping University Electronic Press
Linköping, Sweden

<http://www.ep.liu.se/ea/cis/2001/007/>

*Published on October 15, 2001 by
Linköping University Electronic Press
581 83 Linköping, Sweden*

**Linköping Electronic Articles in
Computer and Information Science**

ISSN 1401-9841

Series editor: Erik Sandewall

©2001 Jacques Cohen

Typeset by the author using L^AT_EX

Formatted using étendu style

Recommended citation:

*<Author>. <Title>. Linköping Electronic Articles in
Computer and Information Science, Vol. 3(2001): nr 7.
<http://www.ep.liu.se/ea/cis/2001/007/>. October 15, 2001.*

This URL will also contain a link to the author's home page.

*The publishers will keep this article on-line on the Internet
(or its possible replacement network in the future)
for a period of 25 years from the date of publication,
barring exceptional circumstances as described separately.*

*The on-line availability of the article implies
a permanent permission for anyone to read the article on-line,
to print out single copies of it, and to use it unchanged
for any non-commercial research and educational purpose,
including making copies for classroom use.*

*This permission can not be revoked by subsequent
transfers of copyright. All other uses of the article are
conditional on the consent of the copyright owner.*

*The publication of the article on the date stated above
included also the production of a limited number of copies
on paper, which were archived in Swedish university libraries
like all other written works published in Sweden.
The publisher has taken technical and administrative measures
to assure that the on-line version of the article will be
permanently accessible using the URL stated above,
unchanged, and permanently equal to the archived printed copies
at least until the expiration of the publication period.*

*For additional information about the Linköping University
Electronic Press and its procedures for publication and for
assurance of document integrity, please refer to
its WWW home page: <http://www.ep.liu.se/>
or by conventional mail to the address stated above.*

Abstract

We investigate a few approaches that have been considered in the simulation and modeling of networks describing cell behavior. By simulation it is meant the direct problem of determining cell behavior when given a graph (network) specifying the interaction among genes. By cell behavior we mean determining the amount of byproducts (mRNA or protein) that each gene generates with time as it interacts with other genes. We refer to modeling as the inverse problem namely, inferring the network graph when given the data describing the cell's behavior. The modeling problem has acquired significant importance in view of the present high volume of cell data available from micro-array experiments. The emphasis of the paper is in using the constraint logic programming paradigm to describe the simulation of cell behavior. In that paradigm the same program describes both a problem and its inverse. Basically one defines multi-dimensional regions, transitions (specifying how control is transferred from one region to the other), and trajectories (sequences of transitions describing cell behavior). The paradigm is applied to several approaches that have been proposed to study simulation and modeling. Several logic programs have been developed to prototype those approaches under the same proposed paradigm. They include considering Boolean and discrete domains. In each case the potential of obtaining practical solutions to the inverse problem are discussed. The proposed paradigm is related to machine learning and to the synthesis of finite-state automata.

1 Introduction

Historically, the first attempts to study cell regulation were based on establishing systems of non-linear differential equations that “simulate” the behavior of the variables involved, as a function of time. In that context the variables correspond to genes capable of producing compounds (mRNA and proteins) that in turn influence the production capacity of other genes. This problem has been thoroughly studied in the mathematical theory of dynamic systems.

The advent of micro-arrays urgently motivates the solution of the *inverse* problem: given the curves expressing the capacity of each gene of producing known compounds as time progresses, attempt to infer the original system of equations or equivalently, the graph expressing the interactions among genes. Presently, biologists are capable of generating huge sets of data expressing the amount of production of proteins or mRNA as a function of time [21]. One then wishes to establish the corresponding gene interaction graph. The inverse problem may be called *modeling* in contrast with *simulation*, the term used to describe the direct problem of determining the behavior of the variables in a system of differential equations.

There have been various attempts to provide discrete solutions to the system of differential equations describing cell regulation. By considering Boolean variables, instead of continuous variables, one can simplify remarkably the complexity of the original problem. Several papers are now available that describe solutions of both the direct and inverse problem using the Boolean approach.

An important aspect of the inverse problem is that the data obtained by biologists is not reliable, in the sense that many experimental errors occur when dealing with thousands of genes. It then becomes imperative to take into account a probabilistic approach of the inverse problem. In other words, given imperfect data expressing the behavior of variables with time, can one infer the corresponding gene interaction graph ?

More recently this last problem has been simplified by having biologists conjecture the interaction graph and propose the simpler question: How likely is the hypothesized graph capable of generating the given data set ? (See for example [19])

An initial classification of the existing methods proposed in the literature to solve the direct and inverse problems involves three parameters and thus eight subclasses:

Probabilistic versus Nonprobabilistic
Continuous versus Discrete
Direct versus Inverse

In a recent article de Jong and Page [7] propose a *Direct-Discrete-Non-probabilistic* type of approach; that amounts to a qualitative (*à la* Kuipers [10]) simulation of the systems of differential equations representing a given influence graph. Even though their algorithms are implemented in a standard programming language, the authors make extensive use of constraints in formulating the solution to that problem.

In this paper we first propose the use of constraint logic programming as a paradigm for solving the discrete versions of the cell regulation problem. Basically one considers *regions* in an n -dimensional space, and *transitions* that specify possible paths between the regions. Both regions and transitions are specified by sets of constraints that restrict the values of the

variables being studied. A path linking an initially specified region to subsequent regions via the transitions, expresses the behavior of the system. (That description is inherent to the approach taken by de Jong and Page.)

The inverse problem amounts to determining regions and transitions when given the paths that may be derived from the data expressing the behavior of the variables. From those regions and transitions one should be able to determine the desired gene interaction graph.

In the particular case of Boolean constraints – a special case among the discrete approaches – the determination of the regions and transitions amounts to the problem of synthesizing Boolean circuits from input and output data. The regions correspond to subsets of the input data and the transitions correspond to Boolean formulas that allow transforming the input data into output. Therefore the problem using Boolean variables can be expressed in terms of program synthesis or machine learning. Similarly, it is hoped that the synthesis of constraints from data should play a role in the discrete solution of the inverse problem.

An interesting payoff of the use of constraints is that they may well pave the way for probabilistic solutions to both the direct and indirect versions of the cell regulation problem. The approach reinforces the importance of developing probabilistic machine learning algorithms expressed in logic programming and in constraint logic programming.

1.1 Objectives

The emphasis of this paper is to provide a unified view of the various approaches that have been proposed in the study of regulation. This is achievable by making extensive use of constraints, in the sense of CLP (Constraint Logic Programming). It will be seen that different versions of the direct problem, namely determining gene behavior as a function of time from given labeled graphs, can be described by the same general principles.

It will be shown in the following sections that most of the existing approaches to solve the direct problem are based on the concepts of *regions* (in n -dimensional space), *transitions* specifying how control moves from one a region to another, and *trajectories*, a sequence of consecutive transitions.

In most cases the inverse problem (modeling) can then be described as in logic programming: if the clauses specifying the program contain no impure constructs (i.e., all predicates are “backtrack able” in the sense of Prolog) then one should be able, at least theoretically, to perform inverse operations that, given the gene behavior with time, determine the possible labeled graphs by generate-and-test techniques. This may be not be practically achievable in the case of large problems using current computers, but the stress on defining clearly the operations involved in the solution of a problem has definite benefits. Incidentally, this is not unlike the *generate-and-test* approaches used in machine learning and inductive logic programming. The proposed framework may also lead to future approaches using synthesis of finite state automata.

The structure of the paper is as follows. First we present the original formulation of the continuous approach using systems of non-linear differential equations. We then provide a simple motivating example that illustrates the concepts of regions, transitions, and trajectories. That is followed by presentations of the available discrete formulations using the proposed general framework of constraints. The respective subsections deal with approaches to solve the inverse problem. In the conclusion we make suggestions for extensions and further work.

2 The Direct Continuous Approach

Graphs expressing the interactions among genes are often used as convenient abstractions enabling the formulation of the appropriate systems of differential equations. The nodes of those graphs correspond to the variables (genes) and the directed edges – labeled by *plus* or *minus* signs – express the positive or negative interactions among the variables. In the case of cell regulation thousands of genes may interact with each other.

The continuous approach is based on the theory of dynamic systems. From a given regulation graph with n nodes one can always establish a set of n differential equations expressing the behavior of the system with time. The solution of these equations for a given initial condition yields the curves depicting how the genes are expressed as a function of time. The difficulty with this approach is that highly non-linear functions (sigmoids) are used to define the variation of each gene's expression with time.

One of the first approaches in simulating cell behavior consisted of setting up a system of non-linear differential equations of the form:

$$dx_i/dt = f_i(\mathbf{x}) - \gamma_i x_i, \quad x_i \geq 0$$

where i denotes the i^{th} gene and \mathbf{x} is the vector (x_1, x_2, \dots, x_n)

The term $-\gamma_i x_i$ states that the concentration of the i^{th} product decreases through spontaneous processes like degradation, diffusion, etc.; f_{ij} is the function specifying a combination of sigmoids (highly non-linear) and which describes the interaction between nodes (genes) i and j ; m is a parameter specifying the steepness of the function around θ_{ij} (see Figure 1)

$$f_{ij} = h^+(x_j, \theta_{ij}, m) = \frac{x_j^m}{x_j^m + \theta_{ij}^m}$$

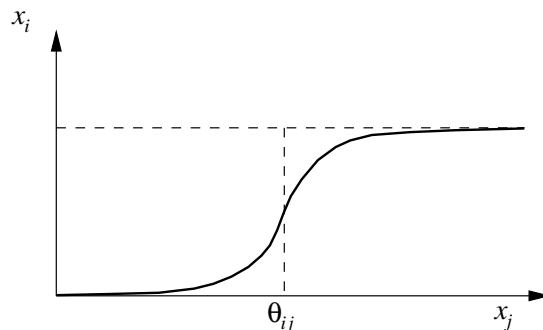


Figure 1: Sigmoid function.

We now open a parenthesis to describe by an example how a set of differential equations can be generated from a labeled graph expressing the interaction between genes. Each gene corresponds to a *node* in the graph, or equivalently, to a *variable* in the system of equations. The simple graph in Figure 2(a) contains three nodes and five edges expressing positive and negative interactions among the genes. The resulting equations are also shown in that figure. Basically, one has to consider different functions h^+

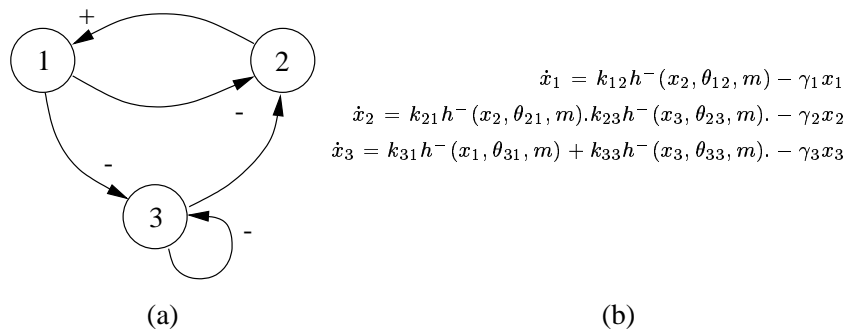


Figure 2: (a) A regulatory network modeled by the 3 equations in (b).

or h^- corresponding to sigmoids that express either a sharp increase (h^+) or a sharp decrease (h^-) of the variable being simulated.

A remark is in order: a simple directed labeled graph does not provide the full information about the nature of the system of equations that can be generated. One can, for example, specify by additional graphical notation that the incoming edges to a node should be considered as being connected by a Boolean *and*, or a Boolean *or*. It will be seen that the Boolean and qualitative discrete type of approaches can incorporate those additional graphical notations. In Figure 2 it is assumed that the two edges entering node 2 are linked by an *and* connector; and the edges entering node 3 are linked by an *or* connector. That explains the product and the sum of the h functions in Figure 2 (b).

It is known from the theory of dynamic systems that the solution of a system with n variables – when given an initial condition specified by the vector \mathbf{x}_0 – can be described by an n -dimensional curve expressing the behavior of the system as a function of time. We call that curve a trajectory. It is known that trajectories may contain loops or may converge to a stationary point.

2.1 The Inverse Problem

It appears that this continuous type of approach makes it extremely difficult to obtain inverse solutions analytically.

2.2 Brief Outline of the Following Sections

In a computer implementation of the direct problem one has to find discrete solutions of the given system of equations. The main approaches that have been utilized in rendering the solution discrete involve either using Booleans or the so-called qualitative methods. Those approaches will be described by examples in later sections of the paper. In the next section we present a motivating example that illustrates the use of constraints in solving a problem that is akin to that of dynamical systems. It will be seen that the constraint formulation in that example can be generalized to cover both the Boolean and qualitative approaches.

3 A Motivating Example

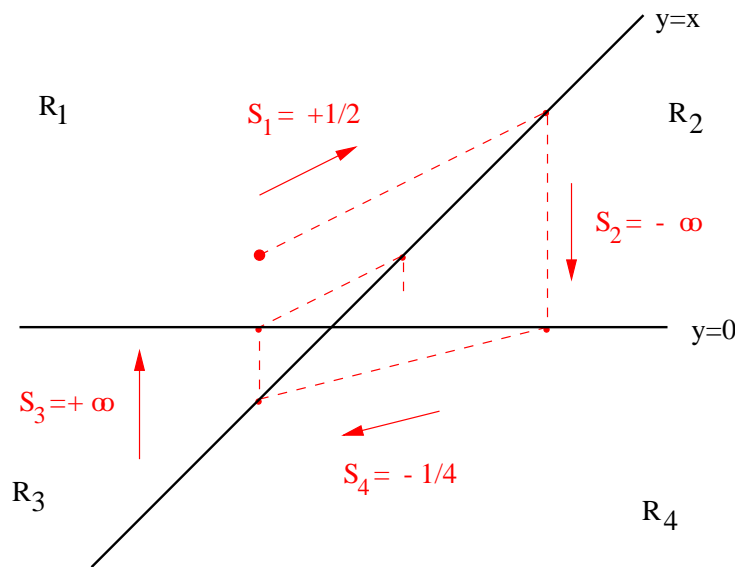


Figure 3: Example of linear constraints and slopes.

Figure 3 presents an intuitive example of regions, transitions, and trajectories in two-dimensional space. In this simplified version of the direct problem one specifies regions by linear constraints. The arrows in that figure specify the transition slopes applicable to each of the four regions. Initially, we deal with deterministic situations, so that there is a single arrow associated to each region.

The direct problem then becomes: given an internal initial point in one of the regions, determine a trajectory, i.e., a sequence of line segments that start at the given point and have the slopes as specified by the arrows. (To avoid ambiguous situations we assume that the arrow associated to a region is not directed towards the intersection of two lines defining that region. One also assumes that a trajectory does not contain any of the lines defining the regions' boundaries.)

Asarin, Maler, and Pnueli [5] have shown that in two-dimensional real space the following property is decidable: given two points P_1 and P_2 one can determine if the trajectory starting from P_1 passes through P_2 . The complexity of the problem is exponential. The authors also proved that that property *is not* decidable in three-dimensional and higher spaces.

Let us assume that regions and their associated slopes are specifiable by constraints; one can then sketch a logic program that determines trajectories that start at a given point. That program is presented in Figure 4.

The predicate `region` specifies a region defined by constraints. The predicate `transition` defines the next region accessible from a given region via a given `slope`. One can then define the trajectory as the sequence of regions that are traversed from a given starting point. One such trajectory is also shown in Figure 3.

In the program of Figure 4, `trajectory` is a predicate that determines a sequence of transitions having n segments, n being specified by a user. Notice that the concept of trajectory is time-related. By considering that


```

region((1,X,Y)):- Y>=0,Y>=X.    region((3,X,Y)):- Y<0,Y>=X.
region((2,X,Y)):- Y>=0,Y<X.    region((4,X,Y)):- Y<0,Y<X.

slope(1,X,Y):- X = 2 * Y.        slope(3,X,Y):- X = 0, Y > 0.
slope(2,X,Y):- X = 0, Y < 0.    slope(4,X,Y):- X = -4 * Y.

trajectory(N,Region,[Region]):- N = 0.
trajectory(N,Start_region,[Start_region|Result]):-
    N > 0,
    region(Start_region),
    transition(Start_region, Next_region),
    region(Next_region),
    N2 = N - 1,
    trajectory(N2,Next_region,Result).

% transition determines the new point of coordinates (X2,Y2)
% which is the intersection of the boundary common to the
% regions N1 and N2 with the straight line of slope(N1,_,_)

transition((N1,X1,Y1),(N2,X2,Y2)):-
    ...

% Note for the particular case in question (Figure 3),
% transition could be written as

transition((R1,X1,Y1),(R2,X2,Y2)):-
    R2 \= R1,
    Delta_X = X2 - X1,
    Delta_Y = Y2 - Y1,
    slope(R1,Delta_X,Delta_Y),
    region((R1,X2,Y2)).

```

Figure 4: Nucleus of a constraint logic program simulating the example in Figure 3.

each transition takes a unit time to be simulated, one can plot a curve expressing how a property of the regions varies with time, given an initial starting point. From now on we may refer to a trajectory as $T(t)$, a function of the variable t expressing time.

3.1 The Inverse Problem

The inverse problem can be stated as follows: given a set of sequences of regions traversed from given initial points determine the constraints that constitute the regions and slopes. In other words, given $T(t)$ determine the corresponding regions and slopes. This is a typical problem in machine learning and a challenge is made to practitioners of inductive logic programming to attempt to solve the problem of Figure 3 in the two-dimensional space using linear constraints. Notice that even though the direct problem is linear, the inverse one *is not*, since the coefficients of the linear equations

and inequations are unknown.

It appears that, in the two-dimensional case, a trajectory can also be expressed by a regular expression in the vocabulary of regions. (An argument to that effect is that the number of regions is finite and that the sequence of regions in a trajectory eventually repeats itself; see Asarin, Maler, and Pnueli [5] for a proof in the case of 2D).

This suggests that the direct problem in 2D consists of finding the regular expressions describing the trajectories when given a set of initial points, the constraints specifying the regions and the respective slopes. The inverse problem would consist of: given the initial states and the regular expressions describing an experiment, determine the constraints and slopes that would yield the given data.

One can also express a sequence of transitions by a directed graph: its nodes represent the regions and the edges any possible transitions. This description corresponds to a finite-state transition graph having unlabeled edges.

A non-deterministic logic program that would try all possible combinations of possible slopes could simulate the non-deterministic case of regions having multiple slopes. This would likely result in a highly combinatorial program with exponential complexity. In that case the finite-state machine describing all possible transitions would be non-deterministic since there may well be two edges emanating from a state representing a region.

Actually one could think of a brute force approach to solving the inverse problem using interval constraints as proposed by Older and Velino [12]. In this case each variable would be initially defined as being in the interval $[-\infty, +\infty]$. The constraints would be used to attempt to narrow down the intervals. Note that interval constraints can be conveniently used to solve non-linear systems of equations, which is the case of the inverse problem.

A more restricted version of the inverse problem would be: given all regions and slopes except that of a single region, determine its value so that the resulting data has a given configuration. In that case one would use interval constraints non-deterministically by splitting the interval being considered for the unknown slope.

4 The Direct Boolean Approach

Let us now consider the case of the Boolean domain with n variables. An n -digit binary number represents each region. A “slope” in the example of Figure 5 can be viewed as a Boolean formula that, when applied to a given binary number, yields another binary number also representing a region. That situation is depicted in Figure 7(a). Notice that essentially the same framework program in Figure 6 is applicable, provided that slopes are defined by Boolean formulas. That example is taken from Glass and Kauffman [9].

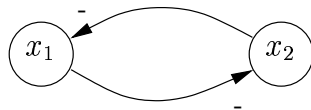


Figure 5: Graph depicting the interaction between two genes.

Actually the Boolean formulas may contain *and* and *or* operators, and can therefore be generated from the directed graph expressing the interactions among genes with additional notations to specify conjunctions and disjunctions, as suggested in the Section 2.

A trajectory, in the sense of the example of Figure 7(a), then becomes a sequence of binary numbers. As in Section 3 individual transitions are assumed to take unit time. It should be obvious that since we are dealing with a finite number (2^n) of regions the determination of all possible trajectories – starting from each possible region – is an exponential problem.

It should also be clear that trajectories could be represented by a finite graph, which may contain loops (including self loops involving a single node). Figure 7(a) represents all the trajectories for the input graph considered in Figure 5. Finally, Figure 7(b),(c) and (d) shows how the two variables change with time when starting from the initial configurations depicted in the figure.

```

region([X1,X2]):- X1=0,X2=1.   region([X1,X2]):- X1=1,X2=1.
region([X1,X2]):- X1=0,X2=0.   region([X1,X2]):- X1=1,X2=0.

slope([X1,X2],[NX1,NX2]):- NX1 = not(X2), NX2 = not(X1).

trajectory(N,Region,[Region]):- N = 0.
trajectory(N,Start_region,[Start_region|Result]):-
    N > 0,
    region(Start_region),
    transition(Start_region, Next_region),
    region(Next_region),
    N2 = N - 1,
    trajectory(N2,Next_region,Result).

% transition determines the region R2 that can be reached
% from region R1.

transition(R1,R2):- slope(R1,R2).

```

Figure 6: Constraint program for the Boolean approach.

4.1 A More Complex Example

It is simple to develop programs similar to the one presented in Figure 6 to deal with more complex graphs. For example, the graph in Figure 8 (a), yields the results presented in (c) and (e). In this particular example, the transitions are given explicitly (*in lieu* of a boolean formula) by the contents of the table in Figure 8 (b).

4.2 The Inverse Boolean Approach

The inverse Boolean domain problem then becomes: given a set of trajectories determine the Boolean formula that represents the transitions between regions. This is a typical problem in circuit design. In that case one deals

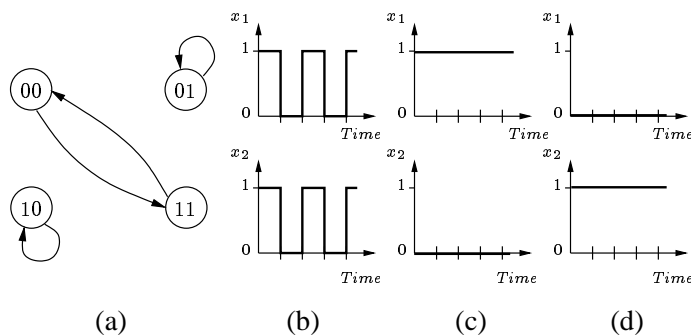


Figure 7: (a) A direct graph depicting all trajectories starting from each region and variations of variables x_1 and x_2 with time, for an initial state $x_1x_2 = 11$ (b), 10 (c), 01 (d) .

with an optimization problem: the determination of a minimal Boolean formula (circuit) capable of yielding a binary number from another one. It is expected that a suitable ILP program could generate the appropriate Boolean formula when given a sequence of pairs of binary numbers.

4.3 The Inverse Boolean Approach Proposed by Somogyi

The members of the group lead by Roland Somogyi have published several interesting papers on the Boolean approach to modeling. Many of those papers are described and referenced in the recent book edited by Bower and Bolouri [1]. In the inverse approach proposed in REVEAL [11], Somogyi introduces the concept of entropy, which is essentially a way to speed up the determination of genes that could possibly interact with other genes.

Given a table such as that depicted in Figure 8 (b), one wishes to determine the corresponding graph, namely that of Figure 8 (a) . An exhaustive generate-and-test approach consists of searching for all Boolean formulas that would express the input-output behavior of the system being considered. By introducing the notion of entropy, Somogyi is able to make that search efficient in cases where a few genes are influenced by a relatively small number of other genes. Essentially, the entropy approach allows one to first consider a single gene that may interact with a given gene; then all pairs of genes, and so forth.

The reference [4] essentially enumerates the possible searches of genes that can influence other genes within the Boolean approach. One can view these efforts as worthy ones in analyzing and improving the efficiency of the inverse Boolean approach. Though Somogyi and his group have been able to deal with modeling problems involving 50 or more genes, much remains to be done by incorporating probabilistic techniques and making it possible to study the interaction of hundreds, and possibly thousands of genes.

5 The Qualitative Approach

With the above background material it becomes easy to explain the approach taken by de Jong and Page [7]. Their so-called qualitative approach

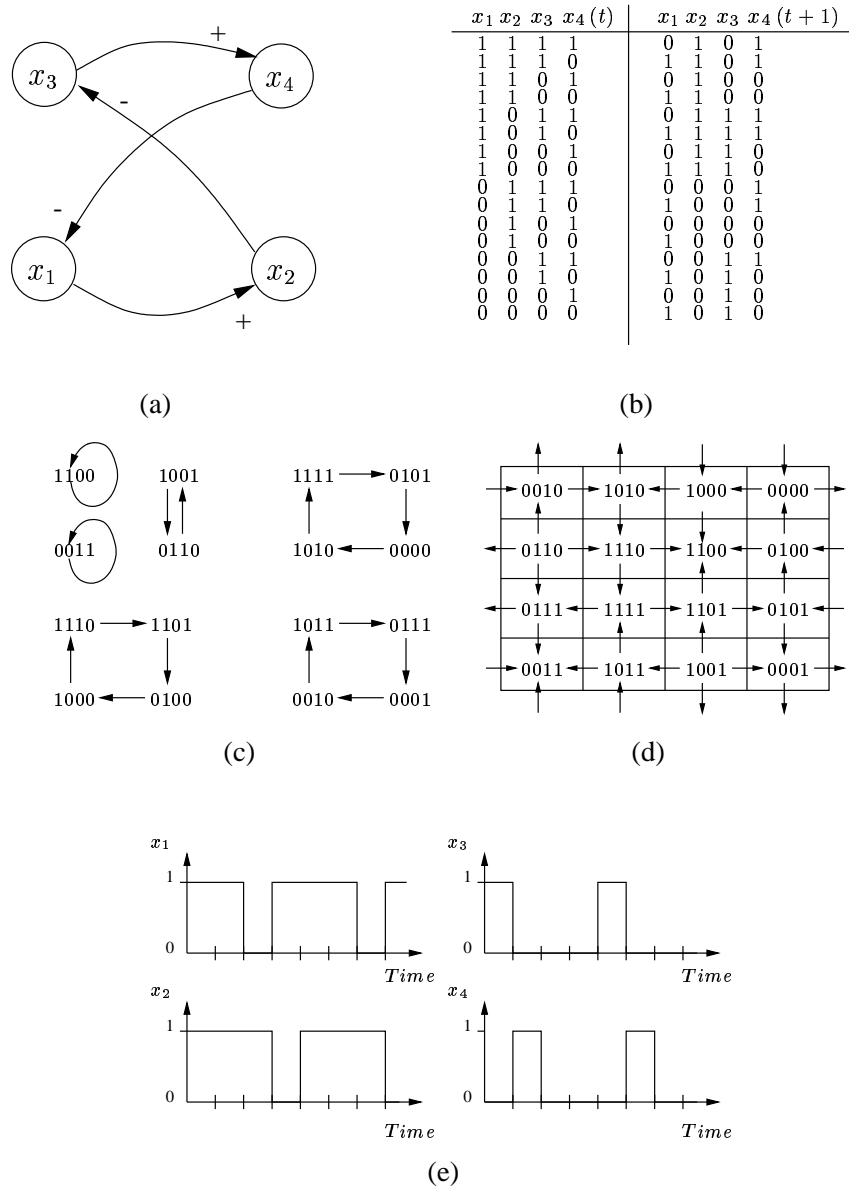


Figure 8: (a) Graph depicting the interaction between four genes and (b) the associated table of transition. (c) A direct graph depicting all trajectories and (d) the mapping of those trajectories on a torus. (e) Curves showing the behavior of the different genes with the initial state 1110.

(in the sense of Kuipers [10]) can be viewed as a generalization of the Boolean approach. Recall that in the latter one assumes that the Boolean formulas representing a transition between regions imply that the n -digit binary numbers involved are changed simultaneously. De Jong and Page assume that each variable representing a gene changes according to a step function simulating a sigmoid function. They also assume that the change of value occurs at a specified position θ (see Figure 1).

In the deJong-Page approach the regions amount to n -dimensional mini-cubes whose edge sizes are determined by the various θ 's. Figure 10 shows a view of the regions involved in the case of 2 variables. As in the example of Figure 10, a trajectory is determined by applying the known slopes (associated to each mini-cube). A slope in that case is defined by qualitative values $q\dot{x}_i$, which may take the values $q\dot{x}_i > 0$, $q\dot{x}_i < 0$, or $q\dot{x}_i \stackrel{\geq}{=} 0$, meaning that the qualitative value qx_i is likely to increase, decrease or to be the same; where qx_i represents a set of constraints specifying a given mini-cube.

```

region((1,X,Y)):- X>=0, X<Theta1, Y=0, Y<Theta2.
region((2,X,Y)):- X>=Theta1, X<max1, Y>=0, Y<Theta2.
region((3,X,Y)):- X>=0, X<Theta1, Y=Theta2, Y<max2.
region((4,X,Y)):- X>=Theta1, X<max1, Y>=Theta2, Y<max2.

slope(1,SX,SY):- SX = 1,   SY = 1.
slope(2,SX,SY):- SX = 0,   SY = 0.
slope(3,SX,SY):- SX = 0,   SY = 0.
slope(4,SX,SY):- SX = -1,  SY = -1.

trajectory(N,Region,[Region]):- N = 0.
trajectory(N,Start_region,[Start_region|Result]):-
    N > 0,
    region(Start_region),
    transition(Start_region, Next_region),
    region(Next_region),
    N2 = N - 1,
    trajectory(N2,Next_region,Result).

transition((R1,X1,Y1),(R2,X2,Y2)):-
    slope(R1,SX1,SY1),
    slope(R2,SX2,SY2),
    successor((R1,X1,Y1),(R2,X2,Y2),SX1,SY1,SX2,SY2).

% the following predicate tests if two regions R1 and R2 are
% adjacent, and that R2 can be reached from R1 according to
% their corresponding slopes.

successor((R1,X1,Y1),(R2,X2,Y2),SX1,SY1,SX2,SY2):-
    ...

```

Figure 9: Constraint program for the qualitative approach.

My student Tri Nguyen-Huu developed a prototype logic program – using the paradigm proposed in this paper – that, given a set of differential equations, discretizes them according to the given step functions (repre-

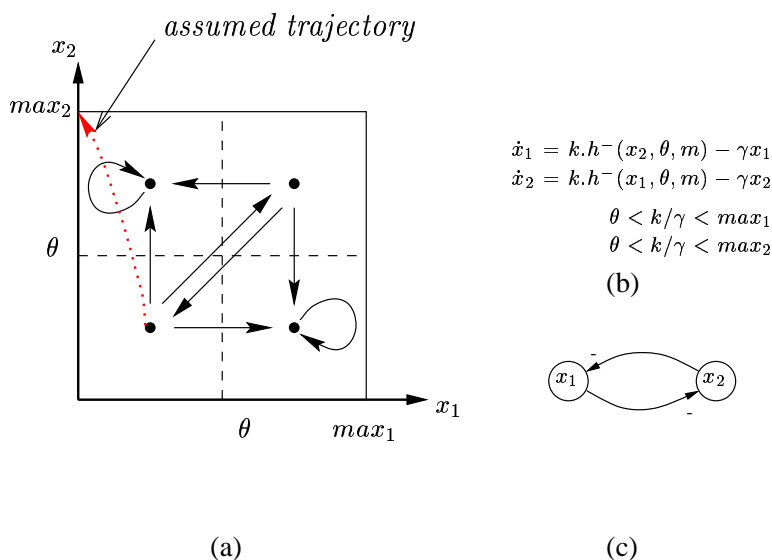


Figure 10: (a) Mini-cubes corresponding to equations in (b) and the interaction graph in (c). The arrows are the slopes associated to each mini-cube.

senting the sigmoids, each with its own θ), and then constructs the set of applicable regions and slopes. In that case each region may be associated with more than one slope. The program then performs a (*don't know*) non-deterministic search capable of determining all trajectories emanating from a given initial point.

The program in Figure 9 can be viewed as having been generated by the Nguyen-Huu's program when given the equations and graph in Figure 10 (b) and (c). The program in Figure 9 has the same general structure as that of the programs in Figures 4 and 6 used to describe the motivating example and the Boolean approach. Assuming the trajectory depicted by the dashed line in Figure 10 (a), one would obtain the curves in Figure 11.

Nguyen-Huu's program is presently non-invertible but we believe that it could be perfected to have that property. Our current goal is to attempt to use it as a tool for helping biologists check if a given set of data is compatible with a hypothesized labeled graph. As a matter of fact this is the approach presently used by de Jong and Page as well as other authors (Gifford [19]).

5.1 The Inverse Qualitative Approach

Consider the differential equation for a variable x :

$$\dot{x} = k - \gamma x$$

A discrete solution means that

$$\frac{x(t+\delta t) - x(t)}{\delta t} = k - \gamma x$$

$$k = \frac{1}{\delta t} x(t + \delta t) + \left(\gamma - \frac{1}{\delta t}\right) x(t)$$

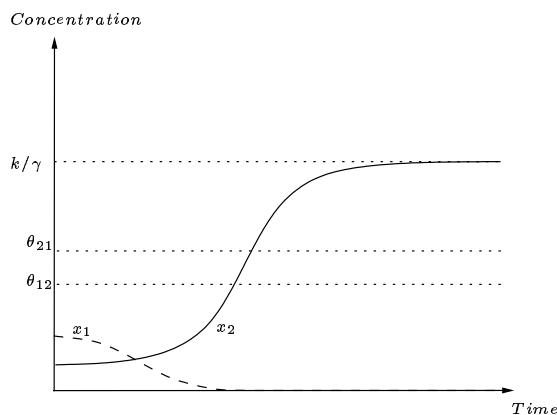


Figure 11: Possible curves obtained for Figure 10.

We first remark the importance of the following property of the solution of the direct qualitative problem:

If one computes the above values of k for a trajectory, then each minicube contains points exhibiting the same value of k .

This is a key property in outlining a possible solution of the inverse qualitative problem. We start by estimating the maximum number of different values of k we can obtain assuming that the deJong-Page approach was used to solve the direct problem. Obviously, one also knows are many genes are involved, say n .

Let us consider the case of two genes and therefore two dimensions. Assuming that the direct qualitative approach has been used, then there are $2^2 = 4$ possible regions in the solution space and there are therefore 4 different values of k . This means that we can label the points in a trajectory constituting data available in the indirect qualitative approach according to those 4 values.

Consider the first gene, and determine the corresponding θ 's so that each of the 4 regions contains the same value of k . That enables the determination of 2 values for θ . Call M_1 the defined regions corresponding to the θ 's associated with the first gene. One then performs a similar operation considering the second gene and thus determines M_2 . The superposition of M_1 and M_2 yields the desired results. This is so because the direct application of the qualitative approach would result in all the regions defined in the superposition of M_1 and M_2 .

In the general case of r edges in the input graph converging to a given node i , one would have 2^r possible values for k . Then one would proceed as above by determining the set of regions M enabling the determination of θ 's for each gene. The superposition of all such M 's yields the desired results.

This is obviously a highly complex set of computations but they outline what is needed to solve the inverse problem using the qualitative approach. This suggests that there may well be approximate less complex solutions to the problem.

6 Comparison Between the Boolean and Qualitative Approaches

A few remarks are in order when comparing the Boolean approach with the qualitative approach. In the Boolean approach there are no specified θ 's defining the time at which there is a change in a Boolean value. In contrast in the qualitative approach time intervals are increased according to the various values of θ that define the faces of the mini-cubes.

A trajectory, in the sense taken in this paper, consists of a traversal in the space that constitutes the minicubes. Trajectories are governed by the current minicube and the possible slopes that it specifies.

Let R_i be the i^{th} region and $Neighbor(R_i)$ be the set of its neighbors. In the qualitative approach the concept of neighboring minicubes is well defined and transitions in a trajectory always take place from neighbor to neighbor.

This may not be the case in the Boolean approach, since one has to clarify the notion of neighborhood. In reference [9] a region is considered as neighboring another one when their binary representations differ by a single digit. This amounts to defining neighborhood in the spatial sense of a torus. Two regions are neighbors if they are adjacent to each other in the torus. Figure 8 (d) depicts the adjacency of the various regions corresponding to the example provided in Figure 8.

7 The General Paradigm Expressed in a Constraint Language

The examples in Sections 3,4 and 5 allow us to generalize the proposed approach using a constraint logic program description. This is depicted in Figure 12.

$$\begin{aligned}
 ®ion_i(\mathbf{x}) \leftarrow region_constraint_i(\mathbf{x}). \\
 &slope_i(\mathbf{x}, \mathbf{x}') \leftarrow slope_constraint_i(\mathbf{x}, \mathbf{x}'). \\
 &transition(region_i(\mathbf{x}), region_j(\mathbf{x}')) \leftarrow slope_i(\mathbf{x}, \mathbf{x}'), \\
 &trajectory(n, region_i(\mathbf{x}), list(region_i(\mathbf{x}))) \leftarrow n = 0. \\
 &trajectory(n, region_i(\mathbf{x}), list(region_j(\mathbf{x}'), result)) \leftarrow \\
 &\quad n > 0, \\
 &\quad transition(region_i(\mathbf{x}), region_j(\mathbf{x}')), \\
 &\quad trajectory(n - 1, region_i(\mathbf{x}'), list((region_j(\mathbf{x}'), result))).
 \end{aligned}$$

Figure 12: General structure of a constraint program traversing regions using transitions (\mathbf{x} and \mathbf{x}' are vectors, i, j and m are integers).

It seems that it would not be extremely difficult to generalize the direct qualitative approach to consider other shapes of sigmoids that would be approximated by piecewise linear segments. In that case one would have to generate a set of regions, slopes and transitions that would constitute the

foundation for the general program computing trajectories as described in Figure 12.

It should be noticed that the programs previously described in this work essentially simulate a finite-state-like automaton. It is convenient to select the set of regions as its terminal vocabulary, as well as its the set of states. One basically has to determine the constraints specifying regions (corresponding to states) and slopes (corresponding to the transitions).

One can then view the inverse problem as the one of synthesizing finite state automata from given data. That synthesis has been used quite often in the past and it represents an alternative to machine learning *à la* ILP. There are Web sites that offer programs capable of performing the synthesis of automata when given thousands of strings [20].

A few words about those programs are in order. First it is important to sketch the functioning of those synthesizers. The given input strings are first represented in the form of a tree that constitutes the automaton accepting exactly the set of input strings, and no other strings. The synthesizer then attempts to merge pairs, then triples, etc., of nodes in the tree therefore introducing loops that are essential for obtaining more general automata. In so doing it is of paramount importance to have a good set of negative examples, otherwise the generated automata is useless since it may consist of a single state. The automata synthesizer approach should be given consideration in future work.

8 Final Remarks

It has been shown in the previous sections that the concepts of regions, transitions, and trajectories were applicable to all the versions of the cell regulation problem considered in this paper. Starting from the original formulation of the direct continuous problem using systems of differential equations, it was shown that the Boolean and qualitative approaches could be described using the constraint logic programming (CLP) paradigm for expressing the relationships between regions, transitions and trajectories.

Under the proposed paradigm both the direct and inverse problems are describable by the same program. We posit that this conceptual development is important because: (1) it presents a unified view of the regulation problem and its possible solutions, and (2) it enables to develop efficient versions of the inverse problem (modeling).

The second point deserves some clarification. Let us consider as an example the case of the Fibonacci function defined by the CLP program for `Fibb(N, Result)`. When solving the inverse problem, i.e., given `Result`, determine `N`, one essentially solves a large system of linear equations. A careful examination of how those equations are solved by an inverse evaluation could lead to more efficient computations.

It was seen in Section 3 that the inverse problem of determining the regions from trajectories is non-linear even though the direct problem is defined by linear inequalities. This suggests the use of interval constraints for solving the inverse problem.

A study of the available literature reveals that clustering plays a key role in the development of solutions to the inverse regulation problem. This was evidenced in the case of the remarks made in Section 5.1 about the inverse qualitative problem. Support Vector Machines (SVM) are related to clustering and have been used in attempts to solve the modeling problem (see [6] and [15]).

The experience gained from the present approach motivates us to extend it to consider the probabilistic case. That is undoubtedly the type of approach needed to solve problems involving a very large number of genes whose regulation data contains experimental errors.

9 Acknowledgements

I wish to acknowledge Tri Nguyen-Huu for all his help in preparing this article. Tri prototyped in Prolog and in CLP the various programs described in this work. The availability of those languages enabled us to develop the programs in a matter of a few weeks.

References

General References (Books, Monographs, Dissertations)

- [1] Bower, J.M., and Bolouri, H., editors, 2001. *Computational Modeling of Genetic and Biochemical Networks*, MIT Press.
- [2] D'haeseleer, P. 2000. Reconstructing Gene Networks from Large Scale Gene Expression Data, Ph.D. dissertation, University of New Mexico, 2000, available at <http://www.cs.unm.edu/~patrik>.
- [3] de Jong, H., 2000. Modeling and Simulation of Genetic Regulatory Systems: A Literature Review, Report 4032, INRIA Alpes.

Articles

- [4] Akutsu, T., Miyano, S., and Kuhara, S. 2000. Algorithms for Identifying Boolean Networks and Related Biological Networks Based on Matrix Multiplication and Fingerprint Function, *Proceedings of the Fourth Annual International Conference on Computational Molecular Biology, Recomb 2000*, Tokyo, ACM Press.
- [5] Asarin, A., Maler, O., and Pnueli, A. 1995. Reachability Analysis of Dynamical Systems having Piecewise Constant Derivatives, *Theoretical Computer Science* vol. 138, pp. 35-65. available at <http://www-verimag.imag.fr//PEOPLE/0ded.Maler>
- [6] Brown M. P. S., *et al*, 2000. Knowledge-based Analysis of Microarray Gene Expression Data Using Support Vector Machines, *Proceedings of the National Academy of Sciences*. vol. 97(1), pp. 262-267.
- [7] de Jong, H., and Page, M. 2000. Qualitative Simulation of Large and Complex Genetic Regulatory Systems, *notes of the 14th Int. Workshop on Qualitative Reasoning*, Morelia, Mexico.
- [8] Friedman, N., Linial, M., Nachman, I., and Pe'er, D. 2000. Using Bayesian Networks to Analyze Expression Data, *Proceedings of the Fourth Annual International Conference on Computational Molecular Biology, Recomb 2000*, Tokyo, ACM Press.
- [9] Glass, L., and Kauffman, S. 1973. The logical analysis of continuous non-linear biochemical control networks. *J. Theor. Biol.* vol. 39, pp. 103-129.

- [10] Kuipers, B. 1994. *Qualitative Reasoning*. MIT Press.
- [11] Liang, S., Fuhrman, S., Somogyi, R. 1998. REVEAL, a General Reverse Engineering Algorithm for Inference of Genetic Network Architectures. *Pacific Symposium on Biocomputing'98* vol. 3, pp. 18-29.
- [12] Older, W., and Vellino, A., 1993. Constraint Arithmetic on Real Intervals. In F. Benhamou and A. Colmerauer, editors, *Constraint Logic Programming: Selected Research*, MIT Press.
- [13] Thieffry D., Thomas R., 1997. Qualitative analysis of gene networks. *Proceedings of the Pacific Symposium on Biocomputing'98*, pp. 77-88.

Web sites

- [14] URL for the Pacific Symposium on Biocomputing
<http://www-smi.stanford.edu/projects/helix/psb98/>
- [15] Support Vector Machines (SVM) URL for basic references
<http://www.eleceng.adelaide.edu.au/Personal/hgchew/svm.html>
- [16] Haussler M., Brown, P., et al, Knowledge-based Analysis of Microarray Gene Expression Data Using Support Vector Machines, available at
http://www.support-vector.net/svm_bbl.html
- [17] Clustering:
<http://genome-www.stanford.edu/~sherlock/tutorial.html>
(tutorial)
http://cne.gmu.edu/modules/dau/stat/clustgalgs/clustgalgs_frm.html (algorithms)
<http://www-genome.wi.mit.edu/MPR/GeneCluster/GeneCluster2.html> (program).
- [18] Tutorial on Bayesian Networks
<http://www.cs.berkeley.edu/~murphyk/Bayes/bayes.html>
- [19] Gifford, D., URL for publications and courses
<http://tesla.lcs.mit.edu/genechips>
and <http://tesla.lcs.mit.edu/6892>
- [20] Synthesis of Finite-State Automata
<http://abbadingo.cs.unm.edu>
and <http://abbadingo.cs.unm.edu/dfa.html>
- [21] Bibliography on Micro Arrays in Biology
<http://linkage.rockefeller.edu/wli/microarray/index.html>