



Brief paper

Robust dynamical network structure reconstruction[☆]Ye Yuan^a, Guy-Bart Stan^b, Sean Warnick^c, Jorge Goncalves^{a,*}^a Control Group, Department of Engineering, University of Cambridge, United Kingdom^b Centre for Synthetic Biology and Innovation, Department of Bioengineering, Imperial College London, United Kingdom^c Information and Decision Algorithms Laboratories, Computer Science Department, Brigham Young University, United States

ARTICLE INFO

Article history:

Received 8 February 2010

Received in revised form

5 February 2011

Accepted 18 February 2011

Available online 9 April 2011

Keywords:

Robust network reconstruction

Noise and unmodelled dynamics

Systems biology

ABSTRACT

This paper addresses the problem of network reconstruction from data. Previous work identified necessary and sufficient conditions for network reconstruction of LTI systems, assuming perfect measurements (no noise) and perfect system identification. This paper assumes that the conditions for network reconstruction have been met but here we additionally take into account noise and unmodelled dynamics (including nonlinearities). In order to identify the network structure that generated the data, we compute the smallest distances between the measured data and the data that would have been generated by particular network structures. We conclude with biologically inspired network reconstruction examples which include noise and nonlinearities.

© 2011 Elsevier Ltd. All rights reserved.

1. Introduction

One of the fundamental interests in systems biology is the discovery of the specific biochemical mechanisms that explain the observed behaviour of a particular biological system. In particular, we consider the problem of reconstructing the network structure from input and partially measured output data of a dynamical system, and in turn uncovering the underlying mechanisms responsible for the observed behaviour. The biological network reconstruction problem challenges come from the necessity to deal with noisy and partial measurements (in particular, the number of hidden/unobservable nodes and their position in the network is unknown) taken from a nonlinear and stochastic dynamical network.

There are several tools in the literature to infer causal network structures. These tools are mainly rooted in three fields: Bayesian inference (Dojer, Gambin, Mizera, Wilczynski, & Tiuryn, 2006; Yu, Smith, Wang, Hartemink, & Jarvis, 2004), information theory (ARACNe Basso et al., 2005; Butte & Kohane, 2000; Faith et al., 2007) and ODE methods (inferelator Bansal & di Bernardo, 2007; Bonneau et al., 2006; di Bernardo et al., 2005; Gardner, Bernardo, Lorenz, & Collins, 2003; Sontag, 2008). Details on these and other

methods can be found in several reviews of the field such as Bansal, Belcastro, Ambesi-Impiombato, and Bernardo (2007), Cantone et al. (2009), De Smet and Marchal (2010), Hecker, Lambeck, Toepfer, Someren, and Guthke (2009) and Koyuturk (2010). The vast majority of network reconstruction methods produce estimates of network structure regardless of the informativity of the underlying data. In particular, most methods produce estimates of network structure even in cases with data from only a few experiments. Such data may not contain enough information to enable the accurate reconstruction of the actual network, thus the obtained network estimates can be arbitrarily different from the true network structure (Cantone et al., 2009). To compensate for the lack of information in data, most methods have heuristics that try to “guess” at the remaining information, either by specifying prior distributions or by appealing to a priori beliefs about the nature of real biological networks, such as looking for the sparsest network. Nevertheless, these heuristics bias the results and lead to incorrect estimates of the network structure.

In contrast, our approach has been to identify the conditions when data is sufficiently informative to enable accurate network reconstruction. The results indicate that even in an ideal situation, when the underlying network is linear and time-invariant (LTI) and the measurements are noise-free, network reconstruction is impossible without additional information (Goncalves & Warnick, 2008). Surprisingly, this information gap is not due to a lack of data, or a deficiency in the number of experiments, but rather it occurs because system states are only partially observed; the information gap is present in all data sets except those that satisfy certain experimental conditions. Our analysis identified a particular experimental protocol that satisfies these necessary conditions to ensure that data will be sufficiently informative to enable network reconstruction. This protocol suggests the

[☆] This work was supported in part by EPSRC grant numbers EP/G066477/1 and EP/I029753/1, AFRL FA8750-09-2-0219 and Microsoft Research through the Ph.D. Scholarship Program. The material in this paper was partially presented at the 49th IEEE Conference on Decision and Control, December 15–17, 2010, Atlanta, Georgia, USA. This paper was recommended for publication in revised form by Associate Editor Elling Jacobsen, under the direction of Guest Editor Francis J. Doyle III.

* Corresponding author. Tel.: +44 1223 3 32770; fax: +44 1223 3 32662.

E-mail addresses: yy311@eng.cam.ac.uk (Y. Yuan), g.stan@imperial.ac.uk (G.-B. Stan), sean@cs.byu.edu (S. Warnick), jmg77@eng.cam.ac.uk (J. Goncalves).

following.

1. A network composed of p measured species demands p experiments.
2. Each experiment requires a distinct input that independently controls a measured species, i.e. experimental input i must affect measured species i and no other measured species except, possibly, indirectly through measured species i .

If data acquisition experiments are not performed in this (or an equivalent) way, the network cannot be reconstructed. Moreover, the resulting information gap is catastrophic, meaning that any internal network structure explains the data equally well (i.e. fully decoupled, fully connected, and everything in between). On the other hand, if some information about the network is available a priori, as is usually the case, then these conditions can be relaxed as explained in Gonçalves and Warnick (2008).

The work in Gonçalves and Warnick (2008), however, did not take into account the realistic scenario that typically systems are nonlinear and data are noisy. This paper extends and details earlier results in Gonçalves and Warnick (2009) by developing an effective method to reconstruct networks in the presence of noise and nonlinearities, assuming that the conditions for network reconstruction presented above in (1) and (2) have been met. Steady-state (resp. time-series) data can be used to reconstruct the Boolean (resp. dynamical) network structure of the system.

The paper is organised as follows. After a motivating example showing that input–output data alone does not enable network reconstruction, Section 2 reviews dynamical structure functions and gives fundamental results concerning their usefulness in the network reconstruction problem. Section 3 presents the main results of the paper regarding robust network reconstruction from input–output data subject to noise and nonlinearities. Finally, we conclude the paper with biologically inspired examples in Section 4.

Notation. For a matrix $A \in \mathbb{C}^{M \times N}$, $A_{ij} \in \mathbb{C}$ denotes the element in the i th row and j th column while $A_j \in \mathbb{C}^{M \times 1}$ denotes its j th column. For a column vector α , $\alpha[i]$ denotes its i th element. We define $e_r^T = [0, \dots, 0, 1_r^{\text{th}}, 0, \dots, 0] \in \mathbb{R}^{1 \times N}$. I denotes the identity matrix. When it is clear from the context, we omit the explicit dependence of transfer functions on the Laplace variable s , e.g. we write G instead of $G(s)$.

Motivating example. Consider the transfer function

$$G(s) = \frac{1}{s+3} \begin{bmatrix} \frac{1}{s+1} \\ \frac{1}{s+2} \end{bmatrix}$$

obtained from data (partial observations) using system identification tools. For simplicity, assume that $G(s)$ accurately represents the input–output relation of the original system. This transfer function is consistent with two state-space realisations $\dot{x} = Ax + Bu$, $y = Cx$ given by

$$A_1 = \begin{bmatrix} -1 & 0 & 1 \\ 0 & -2 & 1 \\ 0 & 0 & -3 \end{bmatrix}, \quad A_2 = \begin{bmatrix} -2 & -1 & 1 \\ -1 & -3 & 1 \\ 0 & -1 & -1 \end{bmatrix}, \quad (1)$$

$B_1 = B_2 = [0 \ 0 \ 1]^T$, and $C_1 = C_2 = [I \ 0] \in \mathbb{R}^{2 \times 3}$ (i.e., the third state is hidden/non-observable). Note that both realisations are minimal and correspond to very different network structures as seen in Fig. 1. This demonstrates that even in the idealised setting (LTI system, no noise and perfect system identification), network reconstruction in the presence of hidden/unobservable states is not possible without additional information about the system.

2. Dynamical structure functions and network reconstruction

In Gonçalves and Warnick (2008) we introduced the notion of dynamical structure functions and showed how they can be used to obtain necessary and sufficient conditions for network

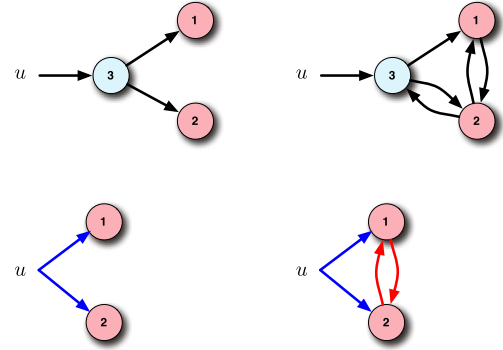


Fig. 1. The same transfer function yields two minimal realisations with very different network structures (left vs. right). Pink nodes are measured (nodes 1 and 2), while blue nodes (here, node 3) represent unmeasured hidden states; the top diagram on either side reveals the complete network structure explicitly showing hidden states, while the lower diagram indicates the corresponding causal structure captured by the dynamical structure function (edges associated with Q are red, while those associated with P are blue). The system on the left is (A_1, B_1, C_1) in (1), and the system on the right is (A_2, B_2, C_2) in (1). Note how completely different the two network structures are (completely decoupled vs. fully connected) even though either realisation would be an equally valid description if all one knew about the system was its transfer function, identified from input–output data. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

reconstruction. For the sake of clarity and completeness, we state these previously obtained results here without proofs. We refer the interested reader to Gonçalves and Warnick (2008) and Yuan, Stan, Warnick, and Gonçalves (2009) for the corresponding proofs.

Consider a nonlinear system $\dot{\tilde{x}} = f(\tilde{x}, \tilde{u}, w_1)$, $\tilde{y} = h(\tilde{x}, w_2)$ with p measured states \tilde{y} , hidden states \tilde{z} (potentially a large number of them), m inputs \tilde{u} , and noises w_1, w_2 . The system is linearised around an equilibrium point (i.e., a point $(\tilde{x}^*, \tilde{u}^*)$ such that $f(\tilde{x}^*, \tilde{u}^*, 0) = 0$), and it is assumed that inputs and noises do not move the states too far from the equilibrium point so that the linearised system is a valid approximation of the original nonlinear system. The linearised system can be written as $\dot{x} = Ax + Bu$, $y = Cx$, where $x = \tilde{x} - \tilde{x}^*$, $u = \tilde{u} - \tilde{u}^*$ and $y = h(\tilde{x}, 0) - h(\tilde{x}^*, 0)$. The transfer function associated with this linearised system is given by $G(s) = C(sI - A)^{-1}B$. When we have partial observations, i.e., when $C = [I \ 0]$, we partition the linearised system equation as follows

$$\begin{bmatrix} \dot{y} \\ \dot{z} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix} + \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} u \quad (2)$$

$$y = [I \ 0] \begin{bmatrix} y \\ z \end{bmatrix}$$

where $x = [y^T \ z^T]^T \in \mathbb{R}^n$, is the full state vector, $y \in \mathbb{R}^p$ is a partial measurement of the state (we assume $p > 1$), z are the $n - p$ “hidden” states, and $u \in \mathbb{R}^m$ is the control input. We restrict our attention to situations where output measurements constitute partial state information, i.e., $p < n$. Taking the Laplace transforms of the signals in (2), solving for Z , and substituting into the Laplace transform of the first equation of yields $sY = WY + VU$, where $W = A_{11} + A_{12}(sI - A_{22})^{-1}A_{21}$ and $V = A_{12}(sI - A_{22})^{-1}B_2 + B_1$. Now, letting D be the matrix composed of the diagonal elements of W , we write $(sI - D)Y = (W - D)Y + VU$. We then obtain $Y = QY + PU$ where

$$Q = (sI - D)^{-1}(W - D) \quad \text{and} \quad P = (sI - D)^{-1}V. \quad (3)$$

Given the system in (2), we define the *dynamical structure function* of the system to be (Q, P) . If all the measured states are removed from the system except for Y_i and Y_j then the transfer function Q_{ij} corresponds to the exact transfer function between Y_j (considered as input) and Y_i (considered as output). The same holds for P in terms of U_j and Y_i .

It can be shown that $G = (I - Q)^{-1}P$ (see Gonçalves & Warnick, 2008). Based on this latter relation, it can be seen that the dynamical structure function of a system contains more information than the transfer function, and less information than the state-space representation (Gonçalves & Warnick, 2008). We can then conclude that, with no other information about the system, neither dynamical nor Boolean reconstruction is possible. Moreover, for any internal structure Q there is a dynamical structure function (Q, P) that is consistent with G , i.e., that satisfies $G = (I - Q)^{-1}P$. In particular, this shows that the use of criteria such as sparsity or decoupledness to guide our selection of a proposal network structure can be misleading. If one were to optimise for decoupledness, for example, a dynamical structure $(0, G)$ could and would always be found, regardless of the true underlying structure. Thus, if we are to use these kinds of criteria, they must be firmly justified a priori.

Proposition 1 (Gonçalves & Warnick, 2008). *Given a $p \times m$ transfer function G , dynamical structure reconstruction is possible from partial structure information if and only if $p - 1$ elements in each column of $(Q, P)^T$ are known that uniquely specify the component of (Q, P) in the nullspace of $[G^T \ I]$.*

The importance of this result is that it identifies exactly what information about a system's structure, beyond knowledge of its transfer function, must be obtained to be able to recover the structure without appeal to a priori assumptions, such as sparsity, or parsimony, etc. This enables the design of experiments targeting precisely the additional information needed for reconstruction. In particular when $p = m$ and G is full rank, we observe that imposing that P is diagonal, i.e., that each input controls a measured state independently, is sufficient for reconstruction.

Corollary 1 (Gonçalves & Warnick, 2008). *If $m = p$, G is full rank, and there is no a priori information about the internal structure of the system, Q , then the dynamical structure can be reconstructed if each input controls a measured state independently, i.e., if, without loss of generality, the inputs can be numbered such that P is diagonal.*

3. Robust network structure reconstruction

In this section, we consider the problem of robustly reconstructing dynamical network structures. Data are obtained from input–output measurements of a noisy nonlinear system. From this type of data we aim to find the internal network structure Q associated with the linearised system (2). To average out the noise, data-collection experiments are repeated N times. For simplicity of exposition, we assume that no a priori information on the internal network structure Q is available. The results still follow if some a priori information about Q is available, and such information can typically be used to relax the experimental protocol according to Proposition 1. Hence, data are collected according to the measurement protocol described in the Introduction.

- (1) The number of distinct data-collection experiments is the same as the number of measured species. This in particular implies that $u(t), y(t) \in \mathbb{R}^p$;
- (2) Each input u_i controls first the measured state y_i so that P is a $p \times p$ diagonal matrix.

In the following two Sections 3.1 and 3.2, we propose two approaches for estimating the dynamical structure function (Q, P) from measured input–output data. The first approach is indirect and involves estimating the transfer function G followed by computing (Q, P) from G . Since some information is lost in the process of estimating G , we consider a second approach where (Q, P) is directly estimated from data (without estimating first G). Concerning the type of input–output data collected, we first consider time-series input–output data and then the special case where only steady-state data are available.

3.1. Dynamical network reconstruction from identified transfer functions

This section describes a method to obtain the dynamical structure function from a stable transfer function G . This transfer function was identified from noisy time-series data using standard system identification tools (Ljung, 1999). According to Corollary 1, if G is full rank there is a unique Q and diagonal P satisfying $(I - Q)G = P$. Since G is an approximation of the actual system, Q and P will typically be mere approximations of the actual dynamical structure function. Moreover, due to noise and unmodelled dynamics, it is likely that Q does not even have the correct Boolean structure. Typically, the internal structure function Q obtained from such a procedure will be fully connected, i.e., all non-diagonal elements of Q will be non-zero.

The main idea to solve the network reconstruction problem from noisy data is the following. For p measured states, Q has $p^2 - p$ unknowns. We want to quantify the smallest distance from G (or directly from the measured data) to all possible Boolean structures (and there are $2^{p^2 - p}$ of them). Some of such distances will be large revealing that the corresponding Boolean structures are unlikely to be the correct structures while other will be small making them candidates for the correct structure.

There are a number of ways to model input–output data with noise and nonlinearities. In order to obtain a convex minimisation problem, we consider the output (could also be input) feedback uncertainty model (Zhou, Doyle, & Glover, 1996). In this framework, the “true” system is given by $(I + \Delta)^{-1}G$, where Δ represents unmodelled dynamics, including nonlinearities, and noise. Based on this choice of dynamic uncertainty, the distance from data to a particular Boolean structure is chosen to be $\|\Delta\|$, in some norm, such that Q obtained from $(I + \Delta)^{-1}G = (I - Q)^{-1}P$ has the desired Boolean structure. We can rewrite the above equation as $\Delta = GP^{-1}(I - Q) - I$. Now, let $X = P^{-1}(I - Q)$. Then the Boolean structure constraint on Q can be reformulated on X , i.e., non-diagonal zero elements in X correspond to those in Q (since $X_{ij} = P_{ii}^{-1}Q_{ij}$ for $i \neq j$).

We can order all Boolean structures from 1 to $2^{p^2 - p}$, and define a set \mathcal{X}_k containing transfer matrices that satisfy the following conditions: (i) for $i \neq j$, $X_{ij}(s) = 0$ if for the considered k th Boolean structure $Q_{ij}(s) = 0$; all other $X_{ij}(s)$ are free variables; (ii) when $i = j$, $X_{ii}(s)$ is a free variable. Hence, the distance from G to a particular Boolean structure can be written as $\alpha_k = \inf_{X \in \mathcal{X}_k} \|GX - I\|^2$, which is a convex minimisation problem with a careful choice of a norm. Next, we show that this problem can be cast as a least squares optimisation problem. If we use the norm defined by $\|\Delta\|^2 = \text{sum of all } \|\Delta_{ij}\|_2^2$, where $\|\cdot\|_2$ stands as the \mathcal{L}_2 -norm over $s = j\omega$, then using the projection theorem (Young, 1988) the problem reduces to

$$\begin{aligned} \alpha_k &= \inf_{X \in \mathcal{X}_k} \|GX - I\|^2 = \inf_{X \in \mathcal{X}_k} \sum_i \|GX_i - e_i\|_2^2 \\ &= \sum_i \inf_{Y_i} \|A_i Y_i - e_i\|_2^2 \\ &= \sum_i \|A_i (A_i^* A_i)^{-1} A_i^* e_i - e_i\|_2^2, \end{aligned}$$

where X_i is the i th column of $X \in \mathcal{X}_k$, Y_i is a column vector composed of the free (i.e., non-zero) elements of X_i , A_i is obtained by deleting the j th columns of G when the corresponding elements X_{ij} are 0 for all j , and $(\cdot)^*$ denotes transpose conjugate. The infimum is achieved by choosing $Y_i = (A_i^* A_i)^{-1} A_i^* e_i$, and $A_i^* A_i$ is always invertible since G is full rank in Corollary 1. If experiments are repeated N times, yielding a transfer function G^i for each experiment, then the above analysis still follows simply by letting $G = [G^1 \ \dots \ G^N]^T$.

3.2. Dynamical network reconstruction directly from time-series data

The previous sections used a two-step approach in which system identification was first used to estimate a transfer function from measured input–output data and then, in a second step, the identified transfer function was used to obtain a dynamical structure function representation of the system which is optimal in terms of a particular metric. This section proposes a method which allows identification of the optimal dynamical structure function representation directly from the measured input–output data. The advantage of this direct network structure reconstruction from data is that no information is lost during the initial transfer function identification stage.

Due to the equivalence between dynamical uncertainty perturbations (Zhou et al., 1996), we are free to choose, without loss of generality, the type of uncertainty perturbation that best suits our needs. For the direct method, instead of a feedback uncertainty as was considered in the previous section, the uncertainty perturbation we are considering here is the additive dynamic uncertainty on the output, i.e., $Y = G_\Delta(U + \Delta)$. In this case, we think about the “distance” in terms of how much we need to change the input (data) to fit a particular Boolean structure. Since $G_\Delta = (I - Q)^{-1}P = X^{-1}$, the equality $Y = G_\Delta(U + \Delta)$ can be written as $\Delta = XY - U$,

where $X \in \mathcal{X}_k$, for some particular Boolean network k . Recall that structural constraints in Q can be imposed directly on X from the equality $X = P^{-1}(I - Q)$. We can therefore use system identification tools for non-causal autoregression models under the structural constraints to identify X (which might be non-causal). In this case, the distance is defined as the maximum likelihood of the estimation problem.

3.3. Penalising connections

The above methodology suffers from a crucial weakness: there are several Boolean structures with distances smaller than or equal to the distance to the “true” network. Indeed, the extra degrees of freedom of the fully connected network allow its corresponding distance α_k to be the smallest of all. This is similar to the noisy data over-fitting problem encountered in system identification where the higher the order of the transfer function, the better the fit. The typical approach in system identification is to penalise higher dimensions and the analogy here is to penalise extra network connections.

If the true network has l non-existent connections (l off-diagonal elements in Q are zero) then there are $2^l - 1$ different Boolean networks that have a smaller or equal distance (due to the additional degrees of freedom provided by the extra connections). When noise is present, then the “true” network will typically have an optimal distance similar to those other l networks. The question of how to find the “true” network thus arises. With repeated experiments, small enough noise (i.e., large enough signal-to-noise ratio) and negligible nonlinearities, the optimal distances of those l networks are comparable, and they are typically much smaller than those of the other networks. To try to reveal the “true” network, one can strike a compromise between network complexity (in terms of the number of connections) and data fitness by penalising extra connections. There are several ways to do this. Here, we consider one of the classical methods known as Akaike’s information criterion (AIC) (Hirotsugu, 1974), or some of its variants such as AICc (which is AIC with a second order correction for small sample sizes), and the Bayesian information criterion (BIC) (Burnham & Anderson, 1998).

The AIC-type approach is a test between models—a tool for model selection. Given a data set, several competing models may

be ranked according to their AIC value, with the one having the lowest AIC being the best. From the AIC value one may typically infer that the best models are in a tie and the rest are far worse, but it would be arbitrary to assign a threshold above which a given model is rejected (Burnham & Anderson, 1998). The AIC value for a particular Boolean network B_k is defined as

$$AIC_k = 2L_k + 2 \ln \alpha_k, \quad (4)$$

where L_k is the number of (non-zero) connections in the Boolean network B_k and α_k is the optimal distance for this Boolean network.

Although finding the optimal distance in the second term of Eq. (4) can be done efficiently, the number of Boolean networks 2^{p^2-p} grows very fast with the number of measured states p . To find the network with the smallest distance it is thus not desirable to compute the optimal distance for each possible Boolean network. Fortunately, there are ways to reduce the number of networks that need to be considered. As we saw in the previous section $\inf_{X \in \mathcal{X}_k} \|GX - I\|^2 = \sum_i \inf_{Y_i} \|A_i Y_i - e_i\|_2^2$ meaning that we can solve each optimisation problem separately. Since each Y_i corresponds to $p - 1$ unknowns in the i th row of Q , this reduces the problem to solving p^{p-1} optimal distances. Finding a polynomial-time algorithm to compute the optimal distance through this method is a subject of current investigation. When it comes to the steady-state case, (Bansal & di Bernardo, 2007) proposed a polynomial-time algorithm to quickly find the ranked solutions at the expense of solution accuracy.

3.4. Boolean network reconstruction from steady-state data

So far we have assumed that time-series data are available. Frequently, however, experimentation costs and limited resources only permit steady-state measurements. In addition, with steady-state measurements it is typically possible to perform a larger number of experiments within the same amount of time, effort and cost. As shown below, most of the connectivity of the network together with the associated steady-state gains (and the associated positive or negative sign) can still be reconstructed from steady-state data. However, no dynamical information will be obtainable. In other words, for most cases we can still recover the Boolean network from steady-state data.

Assume that after some time of maintaining the control input concentrations at a constant value, the measured outputs y have converged to a steady-state value. This is equivalent (if the system is stable or quasi-stable Sontag, 2008) to assuming that we can obtain $G(0)$, i.e., $G(s)$ evaluated at $s = 0$. Now the relationship $(I - Q(s))G(s) = P(s)$ evaluated at $s = 0$ becomes $(I - Q(0))G(0) = P(0)$. From this equation and the knowledge of $G(0)$, all of the results given in Sections 3.1 and 3.2 follow provided that no element of $G(s)$ has a system zero (Zhou et al., 1996) at 0. In that case, a non-zero element in the obtained Boolean network indicates the existence of a causal relationship between the corresponding pair of nodes while a zero element indicates the absence of such relationship.

4. Biologically inspired examples

This section illustrates with two examples the theoretical results presented in the previous section. The corresponding sets of ordinary differential equation describing the dynamics of the considered networks are used to generate noisy data, which are then fed to our reconstruction algorithm in order to assess its ability to recover the correct network structure.

4.1. Single feedback loop

In this first example, we consider the following nonlinear system:

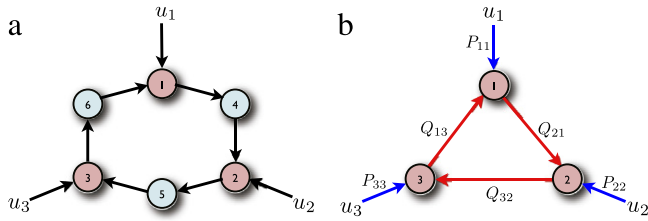


Fig. 2. (a) Complete network with all the states. The red circles represent the measured states (nodes 1, 2 and 3) while the blue circles (nodes 4, 5 and 6) correspond to hidden states. (b) Network of the measured states only. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

$$\dot{y}_1 = -y_1 + \frac{V_{\max}}{K_m + z_3^3} + u_1 \quad (5)$$

$$\dot{y}_2 = -2y_2 + 1.5z_1 + u_2 \quad (6)$$

$$\dot{y}_3 = -1.5y_3 + 0.5z_2 + u_3 \quad (7)$$

$$\dot{z}_1 = 0.8y_1 - 0.5z_1 \quad (8)$$

$$\dot{z}_2 = 1.2y_2 - 0.8z_2 \quad (9)$$

$$\dot{z}_3 = 1.1y_3 - 1.3z_3 \quad (10)$$

where $V_{\max} = 0.5$ and $K_m = 0.1$. Eq. (5) includes a nonlinear function of z_3 known as a Hill equation. It represents a negative regulation of the rate of reaction of y_1 by z_3 . For simplicity, all other terms are linear. In this example, $p = 3$, i.e., there are three measured states (y_1, y_2 and y_3) while the other 3 states are hidden (z_1, z_2 and z_3). The corresponding network is given in Fig. 2(a).

Three experiments were performed. In each experiment, one input was a step while the others were set to zero and data was collected for each of the measured species. The experiments were repeated 3 times to average out the noise. For simplification, in this example, only steady-state data was used. Data was obtained by numerically integrating the differential equations in (5)–(10) and adding independent Gaussian noises. The ratios between standard deviations and means of the steady-state data were within the range [0.35, 1.15], which shows that noise is considerable.

Since the true network has 3 elements in Q equal to zero, there are $2^3 = 8$ networks with a better or equal optimal cost. Computing the corresponding distances and AICc values for all possible Boolean structures between the three measured species, we observed that the distance decreased by an order of magnitude when we arrived at the true network. In addition, AIC, BIC and in particular AICc were able to pick the correct network.

4.2. Chemotaxis in *Rhodobacter sphaeroides*

This section considers the reconstruction of the biochemical network responsible for chemotaxis in *Rhodobacter sphaeroides*. The network is represented in Fig. 3 (see Roberts et al., 2009; Wadhams & Armitage, 2004, for a detailed explanation of this model and its biological interpretation). It involves 10 species dynamically interacting through a complex set of interconnections. To illustrate our method, consider noisy data from 3 species only: Y_3^p , Y_6^p and the “motor” (circled in red in Fig. 3(a)), obtained based on simulations of the nonlinear ordinary differential equation model proposed by Roberts et al. (2009). We follow our prescribed experimental protocol and, for simplification, only steady-state data are used. Relatively large Gaussian noise was added to the collected data to simulate measurement noise in the data set.

Based on the complete network given in Fig. 3(a), the correct network to recover is presented in Fig. 3(b). Computing the corresponding distances and AICc values for all the $2^6 = 64$ possible Boolean networks, we observed that the network with the smallest AICc was not the correct network in Fig. 3(b) as it was missing the Q_{12} link. A closer look at the noisy steady-state data of Y_3^p (from a step input in u_2) revealed an extremely large ratio

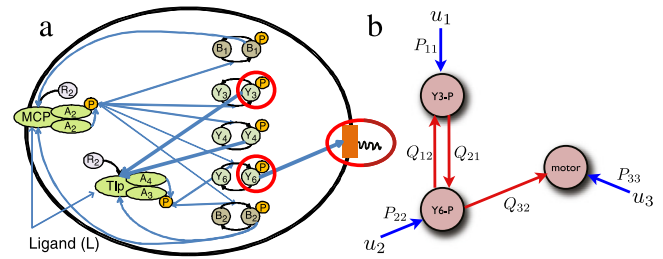


Fig. 3. (a) Network representing the dynamical interaction between the 10 species believed to be responsible for the chemotactic response of *Rhodobacter sphaeroides*. We assume that only species Y_3^p , Y_6^p and “motor” are measured (circled in red). (b) Network connecting the measured states only. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

between its standard deviation and mean value (≈ 200), showing that the noise was completely overpowering the signal. Indeed, it can be shown that Y_6^p has a very small influence on Y_3^p since the pathway from Y_6^p to Y_3^p includes a reversible reaction with a very small rate constant (for detail, see Nordling & Jacobsen, 2009; Yuan, Stan, Warnick, & Goncalves, 2010). The next set of smallest values of AICc consists of 4 networks, including the true one. If necessary, an extra experiment can be performed to further discriminate between these five candidate networks.

5. Conclusion and discussion

This paper proposes a new network reconstruction method in the presence of noise and nonlinearities based on dynamical structure functions. The key idea is to find minimal distances between the existent data and the data required to obtain particular network structures. The method was illustrated with two biologically oriented examples. They showed that even in the presence of nonlinearities and considerable noise network reconstruction was possible. Eventually, when the signal-to-noise ratio was too small, reconstruction was no longer possible, but that is true irrespective of the method used.

Obviously, the method has limitations with respect to nonlinearities. With stronger nonlinear terms eventually the method fails. For example, network reconstruction for oscillatory systems is still an open problem. However, when applied to the reconstruction of various equilibrium point models given in the literature, we observed that reconstruction was always possible when the signal-to-noise ratio of the measured data was not too small (far less than 1).

A final note regarding the application of this methodology to real data. We have looked throughout the literature for real data and none of the available data that we found satisfied the conditions necessary for accurate network reconstruction. Some problems that we observe in the literature include: (1) many publications do not include raw data (they typically only include means and standard deviations), and many authors indicate that they no longer have their data; (2) some microarray data do not include repeats and others were obtained using dual channel microarrays that only give ratios between channels, making it impossible to reliably extract gene expression intensities. One of the most promising papers was Cantone et al. (2009), which followed our experimental protocol, and their raw data is available. We found, however, that the data did not meet the conditions necessary for network reconstruction. In the paper, the authors compared different network reconstruction methods only to find that none of the methods even came close to identifying the true network. Nevertheless, because the authors compared the results to random guessing, they report that “Reverse engineering based on differential equations and Bayesian networks correctly inferred regulatory interactions from experimental data”. We disagree

with their conclusion, and point out that because the data was not sufficiently informative in the first place, such a comparison is not meaningful. To better understand the degree of the lack of information in the data, we considered all 10 subnetworks consisting of 3 nodes. The gap in distances between fully decoupled and fully connected networks ranged from just 2% to a maximum of 70%. This shows that there is not enough information in the data to differentiate between Boolean structures. Note that this data was obtained from over-expression, so based on these results, we hypothesise that over-expressing genes saturate the translation and transcription machinery, making linearisation a poor approximation of the actual system dynamics. Current work is exploring the design of experiments on known systems that (1) satisfy our data-collection protocol to ensure that the resulting data is sufficiently informative for network reconstruction, and (2) facilitate a comparison of various methods so we can better understand how different techniques perform in situations where accurate network reconstruction is, in fact, possible.

References

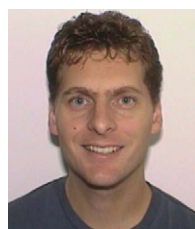
- Bansal, M., Belcastro, V., Ambesi-Impiomato, A., & Bernardo, D. (2007). How to infer gene networks from expression profiles. *Molecular Systems Biology*, 3, 78.
- Bansal, M., & di Bernardo, D. (2007). Inference of gene networks from temporal gene expression profiles. *NET Systems Biology*, 1(5), 306–312.
- Basso, K., Margolin, A., Stolovitzky, G., Klein, U., Dalla-Favera, R., & Califano, A. (2005). Reverse engineering of regulatory networks in human B cells. *Nature Genetics*, 37(4), 382–390.
- Bonneau, R., Reiss, D., Shannon, P., Facciotti, M., Hood, L., Baliga, N., & Thorsson, V. (2006). The inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets de novo. *Genome Biology*, 7, R36.
- Burnham, K., & Anderson, D. (1998). *Model selection and inference—a practical information-theoretic approach*. Springer-Verlag.
- Butte, A., & Kohane, I. (2000). Mutual information relevance networks: functional genomic clustering using pairwise entropy measurements. In *Pac. symp. biocomput.* pp. 418–429.
- Cantone, I., Marucci, L., Iorio, F., Rucci, M., Belcastro, V., Bansal, M., Santini, S., Bernardo, M., Bernardo, D., & Cosma, M. (2009). A yeast synthetic network for in vivo assessment of reverse-engineering and modeling approaches. *Cell*, 137, 172–181.
- De Smet, R., & Marchal, K. (2010). Advantages and limitations of current network inference methods. *Nature Reviews Microbiology*, 8.
- di Bernardo, D., Thompson, M., Gardner, T., Chobot, S., Eastwood, E., Wojtovich, A., Elliott, S., Schaus, S., & Collins, J. (2005). Chemogenomic profiling on a genome-wide scale using reverse-engineering gene networks. *Nature Biotechnology*, 23(3).
- Dojer, N., Gambin, A., Mizera, A., Wilczynski, B., & Tiuryn, J. (2006). Applying dynamic Bayesian networks to perturbed gene expression data. *BMC Bioinformatics*, 7.
- Faith, J., Hayete, B., Thaden, J., Mogno, I., Wierzbowski, J., Cottarel, G., Kasif, S., Collins, J., & Gardner, T. (2007). Large-scale mapping and validation of *Escherichia coli* transcriptional regulation from a compendium of expression profiles. *PLoS Biology*, 5(1).
- Gardner, T., Bernardo, D., Lorenz, D., & Collins, J. (2003). Inferring genetic networks and identifying compound mode of action via expression profiling. *Science*, 301.
- Gonçalves, J., & Warnick, S. (2008). Necessary and sufficient conditions for dynamical structure reconstruction of LTI networks. *IEEE Transactions on Automatic Control*, 53.
- Gonçalves, J., & Warnick, S. (2009). In B. Ingalls, & P. Iglesias (Eds.), *Control theory and systems biology. System theoretic approaches to network reconstruction*. MIT Press.
- Hecker, M., Lambeck, S., Toepfer, S., Someren, E., & Guthke, R. (2009). Gene regulatory network inference: data integration in dynamic models—a review. *BioSystems*.
- Hirotsugu, A. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6), 716–723.
- Koyuturk, M. (2010). Algorithmic and analytical methods in network biology. *WIREs Systems Biology and Medicine*, 2(3), 277–292.
- Ljung, L. (1999). *System identification—theory for the user*. Prentice Hall.
- Nordling, T., & Jacobsen, E. (2009). Interampatternness—a generic property of biochemical networks. *NET Systems Biology*, 3(105).
- Roberts, M., August, E., Hamadeh, A., Maini, P., McSharry, P., Armitage, J., & Papachristodoulou, A. (2009). A model invalidation-based approach for elucidating biological signalling pathways, applied to the chemotaxis pathway in *R. sphaeroides*. *BMC Systems Biology*, 3(105).
- Sontag, E. (2008). Network reconstruction based on steady-state data. *Essays in Biochemistry*, 45, 161–176.
- Wadhams, G., & Armitage, J. (2004). Making sense of it all: bacterial chemotaxis. *Nature Reviews Molecular Cell Biology*, 5, 1024–1037.
- Young, N. (1988). *An introduction to Hilbert space*. Cambridge university press.
- Yu, J., Smith, V., Wang, P., Hartemink, A., & Jarvis, E. (2004). Advances to Bayesian network inference for generating causal networks from observational biological data. *Bioinformatics*, 20(18), 3594–3603.
- Yuan, Y., Stan, G., Warnick, S., & Goncalves, J. (2009). Minimal dynamical structure realisations with application to network reconstruction from data. In *Proceedings of conference on decision and control*.
- Yuan, Y., Stan, G., Warnick, S., & Goncalves, J. (2010). Robust dynamical network reconstruction from noisy data. In *Proceedings of conference on decision and control*.
- Zhou, K., Doyle, J., & Glover, K. (1996). *Robust and optimal control*. Prentice Hall.



Ye Yuan was born on October 1986. He received his B.Eng. degree (Valedictorian) from the Department of Automation, Shanghai Jiao Tong University in 2008, M. Phil. from the Department of Engineering, Cambridge University in 2009. He is currently a second-year Ph.D. student in Control Group, Department of Engineering, University of Cambridge. Ye was a visiting student at University of New Mexico, the Hong Kong University of Science and Technology, and Luxembourg Centre for Systems Biomedicine. He is now a visiting student at CDS, Caltech. His research interest lies in the mathematical control theory with applications to network and biology. He is the recipient of Dorothy Hodgkin Postgraduate Awards (Microsoft Research Ph.D. Scholarship), Cambridge Overseas Scholarship and Henry Lester Scholarship.



Guy-Bart Stan was born in Liège, Belgium, in 1977. He received his Ph.D. degree in Applied Sciences (Analysis and Control of Nonlinear Dynamical Systems) from the University of Liège, Belgium in 2005. In 2005, Dr. Stan worked as a Senior Digital Signal Processing Engineer at Philips Applied Technologies, Leuven, Belgium. From 2006 until 2009, he worked as Research Associate in the Control Group of the Department of Engineering at the University of Cambridge, UK, being supported by an EU-FP6 IEF Marie-Curie Fellowship and the UK EPSRC, successively. In 2008, he was a Visiting Scientist at the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, USA. Since 2009 Dr. Stan is a University Lecturer in Engineering Design for Synthetic Biology Systems in the Department of Bioengineering and the EPSRC-funded Centre for Synthetic Biology and Innovation at Imperial College London. His current research interests include the mathematical modelling, analysis and control of complex biological systems/networks occurring in synthetic biology, systems biology, and technological systems.



Sean Warnick received the B.S.E. degree from Arizona State University in 1993, and the S.M. and Ph.D. degrees in Electrical Engineering and Computer Science from the Massachusetts Institute of Technology in 1995 and 2003, respectively. He attended ASU on scholarship from the Flinn Foundation, graduated summa cum laude, and was named the Outstanding Graduate of the College of Engineering and Applied Sciences. Since 2003 Dr. Warnick has been with the Computer Science Department at Brigham Young University, where he is currently an Associate Professor and Director of the interdisciplinary Information and Decision Algorithms Laboratories. Dr. Warnick has also held visiting appointments at Cambridge University (Summer 2006), the University of Maryland (Summer 2008), and the National Security Agency, where he was named the Distinguished Visiting Professor for three consecutive years (2008–2010). Dr. Warnick's research interests focus on complex networks of uncertain dynamical systems, where he considers issues of representation, identification and reconstruction, control, security and robustness, and verification and validation.



Jorge Goncalves received his Licenciatura (5-year S.B.) degree from the University of Porto, Portugal, and the M.S. and Ph.D. degrees from the Massachusetts Institute of Technology, Cambridge, MA, all in Electrical Engineering and Computer Science, in 1993, 1995, and 2000, respectively. He then held two postdoctoral positions, first at the Massachusetts Institute of Technology for seven months, and from May 2001 to March 2004 at the California Institute of Technology with the Control and Dynamical Systems division. Since April 2004 he has been a Lecturer in the Information Engineering Division of the Department of Engineering at the University of Cambridge. Since 2005 he has been also a Fellow of Pembroke College, Cambridge. He was a visiting researcher at the University of Luxembourg from June to December 2010. Also he is a visiting researcher at California Institute of Technology (January–September 2011). He was the recipient of the Best Student Paper Award at the Automatic Control Conference, Chicago, IL, June 2000.

His research interests include modelling, analysis and control of complex and hybrid systems. In particular, modelling and analysis in systems and synthetic biology, collaborating with biologists in different areas such as circadian rhythms and gene regulatory networks.