



Research

Cite this article: Sachs K, Itani S, Fitzgerald J, Schoeberl B, Nolan GP, Tomlin CJ. 2013 Single timepoint models of dynamic systems.

Interface Focus 3: 20130019.

<http://dx.doi.org/10.1098/rsfs.2013.0019>

One contribution of 11 to a Theme Issue 'Integrated cancer biology models'.

Subject Areas:

computational biology

Keywords:

structure learning, Bayesian networks, perturbations, signalling, networks

Author for correspondence:

C. J. Tomlin

e-mail: tomlin@eecs.berkeley.edu

[†]These authors contributed equally to this study.

Electronic supplementary material is available at <http://dx.doi.org/10.1098/rsfs.2013.0019> or via <http://rsfs.royalsocietypublishing.org>.

Single timepoint models of dynamic systems

K. Sachs^{1,†}, S. Itani^{2,†}, J. Fitzgerald³, B. Schoeberl³, G. P. Nolan¹
and C. J. Tomlin²

¹Department of Microbiology and Immunology, Stanford University School of Medicine, Stanford, CA, USA

²Department of Electrical Engineering and Computer Sciences, University of California at Berkeley, Berkeley, CA, USA

³Merrimack Pharmaceuticals, Cambridge, MA, USA

Many interesting studies aimed at elucidating the connectivity structure of biomolecular pathways make use of abundance measurements, and employ statistical and information theoretic approaches to assess connectivities. These studies often do not address the effects of the dynamics of the underlying biological system, yet dynamics give rise to impactful issues such as timepoint selection and its effect on structure recovery. In this work, we study conditions for reliable retrieval of the connectivity structure of a dynamic system, and the impact of dynamics on structure-learning efforts. We encounter an unexpected problem not previously described in elucidating connectivity structure from dynamic systems, show how this confounds structure learning of the system and discuss possible approaches to overcome the confounding effect. Finally, we test our hypotheses on an accurate dynamic model of the IGF signalling pathway. We use two structure-learning methods at four time points to contrast the performance and robustness of those methods in terms of recovering correct connectivity.

1. Introduction

Learning the structure of biomolecular pathways such as genetic regulatory pathways and signalling pathways is an important task, leading to greater understanding of molecular interactions, improved insight into disease states and increased ability to intervene towards therapeutic ends, for instance, in diseases such as cancer. Careful studies over several decades have revealed many details of biological pathways, leading to revolutionary targeted treatments in some cancer types [1,2], yet due to insufficient detailed information, effective targeted therapies cannot be found for most tumours. An effective automated approach for network elucidation, therefore, can have far-reaching implications. A number of approaches have been employed for the task of structure learning, including probabilistic-based methods such as Bayesian networks (BNs) [3–5], correlation-based methods [6] and mutual information-based methods [7–10]. These methods use *abundance* data, data in which the amount of the biomolecules of interest has been measured. The underlying biological networks are dynamic systems, and some methods such as dynamic Bayesian networks (DBNs) [11] directly address the dynamics of the network. However, data appropriate for learning dynamic models can be difficult to obtain. Even in the data-abundant case of high-throughput single-cell data (as in [3,12,13]), data for learning dynamic models are difficult to obtain. This is because there does not exist a reliable method for connecting cells measured at one timepoint to cells measured at a future timepoint. Real-time live imaging (movies) of cells *does* provide the necessary data; however, it is usually not possible to quantify the relevant variables in live images. Therefore, in many of these studies, the modelling efforts do not address the underlying dynamics; instead, time is treated implicitly, in that the methods tackle the task of learning a structure that summarizes all the influences in the network that occur over a span of time relevant for their context. The models are 'static' (they do not address dynamics)

but the biological systems are dynamic. The methodologies applied do not require data from multiple timepoints; rather, they learn the structure from single timepoint measurements. In this work, we address only this situation, in which appropriate dynamic data cannot be obtained, and the network is learned from measurements taken at a single moment in time from (non-synchronized) individual measurements. We will refer to these single timepoint models as *snapshot* models.

To our knowledge, despite the prevalence of these models, no study to date has addressed the possible impact of dynamics on structure-learning efforts. Studies using snapshot models implicitly assume that statistical and/or information theoretic tools can be used to learn accurate influence connections from the available snapshot data. Is this the case? If so, will any timepoint within a reasonable range yield the same correct influence relationships? These models often assume their data come from *steady state*, but the dynamics of the underlying biological system are arguably still in motion; furthermore, the choice of timepoint at which the system is assumed to have reached steady state may be uncertain, and some systems may not have a well-defined steady state. Additionally, in cases where time course data are available, a dynamic model may not be applied owing to constraints of insufficient data, as is the case for the examples presented below. Can a dynamically faithful snapshot model be learned at each timepoint? Because of the prevalence and usefulness of snapshot models, a close examination of these issues can have important practical impact.

The main focus of our analysis is the model's representation of causal relationships: our models will always reflect statistical relationships, but in order to have impact in disease and biology, we require the models to reflect underlying causality. Several noteworthy studies have addressed the ability of network inference approaches to reflect causality correctly. These include insightful discussions of model interpretability and the impact of steady-state assumptions [14,15]. Still others have discussed the impact of dynamics in *dynamic* models [16–18]. However, none of these addresses the distinct analysis of the impact on dynamics on single timepoint-derived snapshot models.

This study investigates the utility and limitations of snapshot models of dynamic systems by addressing the questions posed above. An empirical study is conducted using synthetic data from a differential equation model of IGF signalling, which closely mimics the behaviour of the biological system, yet provides a reliable ground truth. Our modelling tools are BNs, which are probabilistic models that have been used to learn network structure in biology, and a new tool called generalized Bayesian networks (GBNs), a generalized form we have recently introduced that enables learning of structures with cycles [19,20]. With abundant synthetic data, we can investigate the ability of a snapshot model to uncover underlying influence connections. The synthetic data employed include simulated perturbations in which the activity of a variable is blocked, meaning that the variable abundance is unaffected, but the variable cannot causally affect its usual target(s), mimicking the common class of perturbations available in molecular biology called *small-molecule inhibitors*. Results are reported over several timepoints, and we compare the results of standard BN learning to those of GBNs. We draw conclusions based on our empirical results, shedding light on potential stumbling blocks of structure-learning approaches. A confounding effect inherent to the learning of snapshot

models of dynamic systems is encountered and described in this study; to our knowledge, this is the first time this confounding effect has been described despite the ubiquitous nature of these models. Lastly, we propose methods for overcoming this potential barrier to casual learning.

2. Background and methods

We present background on the modelling methods employed for snapshot modelling, as well as on the source of the synthetic data used in this study.

2.1. Bayesian networks

The first method that we employ in this paper to perform structure learning for snapshot models is BN learning. We note that although we focus on BNs, our results shed light on other techniques which rely on statistical or information theoretic dependencies to elucidate interactions.

BNs [21] represent probabilistic dependence relationships among multiple interacting components (in our case, biomolecules such as mRNA molecules or proteins), illustrating the effects of pathway components upon each other in the form of an influence diagram, or a graph (G) and a joint probability distribution. In the graph, the nodes represent variables (the biomolecules) and the edges represent dependencies (or more precisely, the lack of edges indicates a conditional independence) [21]. For each variable, a conditional probability distribution (CPD) quantitatively describes the form and magnitude of its dependence on its parent(s). Owing to the factorization of the joint probability distribution, the graph must be *acyclic*, meaning that it must not be possible to follow a path from any node back to itself [21].

Since the seminal work by Pe'er and co-workers [4], BNs have been used extensively in biology, to model regulatory pathways both in the genetic [4,22] and in the signalling pathway domain [3,5,23]. The structure-learning task consists of searching the space of possible structures to find the one that best reflects probabilistic relationships in a biological dataset. Under appropriate conditions and assumptions, a causal model interpretation can be applied, indicating that the parent of a variable in the learned graph causally influences the variable's quantity (directly or indirectly) [24,25]. These are known as causal BNs. Within this framework, structure learning can be used to elucidate the structure of interactions in regulatory pathways.

2.2. Generalized Bayesian networks

When building models of pathways, BN models have a number of strong advantages. They are flexible and interpretable, they can handle interactions of arbitrary complexity (given sufficient data), and they can smoothly incorporate both prior knowledge and interventional data in a principled way. However, they have one serious drawback for modelling biological systems: they are unable, as described above, to handle cycles in a static model. Because cycles abound in biological pathways, a BN snapshot model usually cannot hope to capture all influence connections: it will most certainly miss, or incorrectly orient, at least one edge from a directed cycle. Methods which can infer cycles do exist [26,27]; however, they are formulated for linear systems; moreover, they are not as useful in the biomedical context

in which diagnosis and prediction constitute a critical aspect of modelling.

To address this problem and enable the use of BN models in a cyclic domain, we recently introduced GBNs, a generalization of BNs to the cyclic domain [19]. In [19] we also presented an algorithm in which the GBN formalism is used to recover causal structure given interventional (single timepoint) data, in acyclic *or* cyclic domains, and a brief overview and explanation of this GBN structure-learning algorithm is presented below (for details, see [19]). In order to overcome the acyclicity constraint in this study, we employ the GBN algorithm in addition to standard BN structure learning, thus enabling us to assess the success of snapshot models in a cyclic domain.

Algorithm 1: GBN Structure Learning

1. Initialization: a causal BN and intervention set I , where I contains nodes that may be perturbed using small-molecule inhibitors.

2. Probing experiments: collect sets of i.i.d. expression profiles at all nodes, under the cases of applying no intervention, as well as applying each single intervention in I .

3. Detect descendants: based upon response of variables to perturbation, recover descendant information for all nodes in I .

4. Identify cycles: based upon perturbations which affect the abundance of the target variable.

5. BN learning: do Bayesian Network learning with the cycles broken (by interventions on nodes we call ‘cycle breakers’) and integrating the descendant information.

6. Close the cycles: determine the correct edges needed to close the cycles, by detecting the children of the cycle breakers. Any node which remains dependent on its ancestors after conditioning on its parent set will have an edge added from the dependent ancestor to the node, as this persistent dependence indicates the ancestor is actually a parent.

The intuition of this algorithm is as follows. Since BN structure learning does not allow for cycles, our goal is to remove (linearize) cycles, employ standard BN structure learning, then, as a final step, elucidate all edges needed to close existing cycles. Briefly, GBN learning first employs perturbations to detect existing cycles and determine a descendant set for all perturbed nodes, or variables. Next, perturbations are used to break all detected cycles, and the now linearized structure is elucidated using BN learning. Finally, cycles are re-closed by assessing conditional independencies between each node and its non-parent ancestors conditioned on its parents. Ancestors which remain dependent after conditioning on the current parent set are added to the variable’s parent set. In [19], it was proved that given enough data and interventions, this algorithm is guaranteed to recover the causal structure of the underlying system.

2.3. Confounders of structure learning

As we are studying the interplay between causal structure learning and dynamics, it is useful to consider various factors that can affect the results and accuracy of causal structure learning. The main putative confounders of causal structure learning are discussed below.

- *Insufficient data or insufficiently informative data.* In our domain of interest of biomolecular networks, it is often difficult to obtain a dataset with a sufficient number of observations. In this case, true (causal) and spurious edges may have an equal amount of support in the data. When this occurs it is difficult to select high scoring models, as many heterogeneous models will obtain similar scores in the structure-learning effort. Even when a large amount of data are available, it can fail to *inform*. For example, if an important dependence between two variables is not expressed under the experimental conditions employed, it becomes undetectable.
- *Insufficient or imperfect perturbations.* Although it is possible to orient edges in the absence of perturbations, it is straightforward to see that learned structures will often contain unoriented edges in the absence of perturbations. Moreover, perturbations employed in the biological context can have two serious flaws: first they may act ‘off target’, meaning that they perturb variables in addition to or other than the variable they are intended to perturb, and, second, they may fail to effectively perturb their intended target.
- *Cycles.* Due to the acyclic constraint on BNs, cycles present in the underlying system may confound the uncovered causal structure. This can be addressed by employing cyclic methods such as GBNs or other cycle-enabling methods such as structural equation models [26,27].
- *Hidden variables.* Perhaps the greatest hindrance to causal structure learning in our domain is the presence of *hidden variables*, variables which participate in the molecular network and causally interact with the measured variables, yet are not themselves measured, either because their importance is unknown, or more often because bandwidth for measurements is limited, or reagents enabling the detection and quantification of the variables are not available. Hidden variables confound causal learning in two main ways. The first is by intervening between a parent node and its child, meaning that if the edge $x \rightarrow y$ is learned, we must interpret it as *x causally affects y OR x causally affects some (set of) hidden variable(s) which then causally affect y* . This is not very problematic as it merely affects our interpretation of edges, but does not confound the causal structure. More problematic is the second effect, in which hidden variables can impose a statistical dependence on two variables by causally affecting both. This can lead to a non-causal edge learned between any two variables who share a hidden co-parent.

In order to isolate the effects of dynamics, we eliminate *all* the putative confounding effects above (except the cycles, which are handled by the GBM algorithm) by employing synthetic data from a differential equation model, described below.

2.4. Synthetic data and model of insulin-like growth factor signalling

To produce synthetic data, we use a mass action kinetic model describing the dynamics of the insulin-like growth factor (IGF) signalling pathway. IGF signalling is important in normal cell physiology, as well as pathological states such as cancer. A schematic of the IGF signalling pathway is shown in figure 1 and is presented in detail in the appendix. Multiple negative feedback loops had been reported in

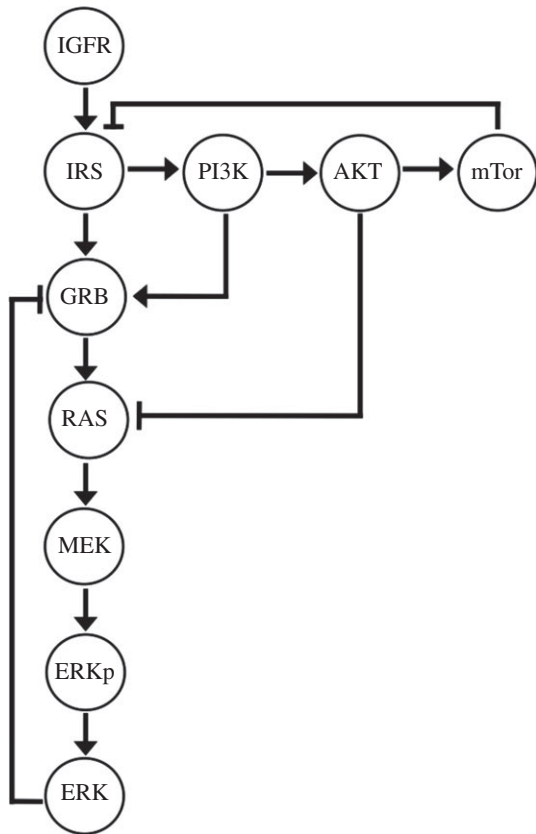


Figure 1. True structure of the underlying dynamic system in IGF signalling. Each node represents the active ‘on state’ of the protein. Perturbations in the form of small-molecule inhibitors are available for MEK, AKT, PI3K, IGFR and mTOR. The simulated data mimics these inhibitors by blocking enzyme activity.

the literature for this receptor/ligand pathway [28,29]. Mass action kinetic equations were used to create the model in Matlab SimBiology v. 2.1. The model contains 60 species, 50 reactions and 37 unique (100 total) parameters. Kinetic parameters reported in the literature were used where available.

There are three directed cycles in the model: $IRS \rightarrow PI3K \rightarrow AKT \rightarrow mTOR \rightarrow IRS$, $GRB2/SOS \rightarrow RasGTP \rightarrow MEK \rightarrow ERKp \rightarrow ERKpp \rightarrow GRB2/SOS$ and $GRB2/SOS \rightarrow RasGTP \rightarrow MEK \rightarrow ERKpp \rightarrow GRB2/SOS$.

The stimulus employed is IGF; in addition, up to five perturbations are employed, at IGFR, MEK, PI3K, AKT and mTOR, corresponding to actual existing small-molecule inhibitors. All of the perturbations are activity inhibitions; that is, they inhibit the protein’s activity, not permitting the targeted protein to phosphorylate other proteins. We generated measurements from four different time points, under 17 total conditions composed of IGF stimulus plus various combinations of inhibitors. For each condition, 1000 unique, randomly selected initial conditions (i.e. molecule concentrations) were employed—the equivalent of collecting 1000 unique cells in a flow cytometry experiment, or performing Western blots on 1000 samples. Simulated ‘measurement noise’ was also added.

The model was created by Jonathan Fitzgerald and colleagues at Merrimack Pharmaceuticals. It is a highly accurate imitation of the true biological system (data not shown), and, accordingly, provides us with synthetic but realistic data, similar to the data one might acquire from a high-throughput measurement technology (as in [3,12,13]). It is a flexible and realistic source of ‘true to life’ synthetic data, which, because

it has a known ground truth model, provides a valuable tool for assessing success of structure-learning efforts.

3. Structure learning

In this section, we study the structure recovered from single-time-point measurements using BNs and GBNs, and contrast it to the known structure of the underlying dynamic system. We study how the resulting structure changes with time and contrasts with the causal dependencies in the dynamic system. While doing that, we study the different structure-learning methods and compare their performance. A total of 17 conditions were used, with no perturbations or single plus multiple perturbations per condition. The models presented are averaged over 20 individual results, edges with confidence more than 0.7 are included. In the following graphs, a dotted edge is a false edge that was predicted and a black edge is a correct edge (possibly inverted). BN structure learning is performed as reported in [3].

3.1. Bayesian networks learning with observational data

We start with observational data in the absence of perturbations. The results are shown in figure 2. Despite the idealized nature of the data from this synthetic system, obtained results deviate sharply from the underlying true structure. Although observational data can only be expected to retrieve a partially directed graph, in this case, many of the expected connections do not appear (approx. half the expected edges are missing in each graph—to simplify the figures, these edges are not indicated), while several non-causal edges do appear. Missed edges, and even non-causal edges, may be attributable to the unmodelled effects of cycles, yet we were nonetheless somewhat surprised by the relative prevalence of non-causal edges. Are all of them caused by the feedback loops? We return to this question below. The graphs vary over time, indicating that in these results, the snapshot model *is* sensitive to the selected timepoint; however, there does not appear to be an indication of ‘early’ versus ‘late’ events (for instance, the influence of IRS on GRB2/SOS is an early step in the model, yet it only appears in timepoint 4), possibly due to the presence of cycles. Therefore, in this example, it is not possible to infer the dynamics of the system from a time course of snapshot models. Although we focused on representing early as well as later timepoints, it is still possible that ‘zooming in’ further on the early timepoints may yield models displaying early and late effects. Additionally, the models are learned from 1000 datapoints—a substantial size dataset by biological standards, but not large for structure learning in general. It is possible that a larger dataset would yield more accurate structures.¹

3.2. Bayesian networks learning with perturbational data

We next consider the results from BN structure learning that includes perturbational data. The results are shown in figure 3. As expected, including perturbational data improves performance substantially. The connectivity structure is recovered reasonably well, with only two or three edges missing in each graph. Still, non-causal edges abound, with about

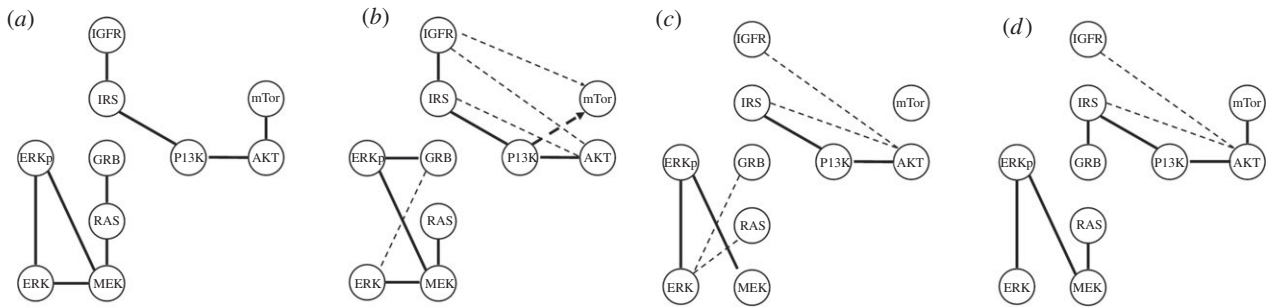


Figure 2. BN structure learning with observational data at different timepoints: (a) $t = 1$, (b) $t = 2$, (c) $t = 3$ and (d) $t = 4$. Solid edges are causally accurate, dotted edges are not. Missed edges are not indicated.

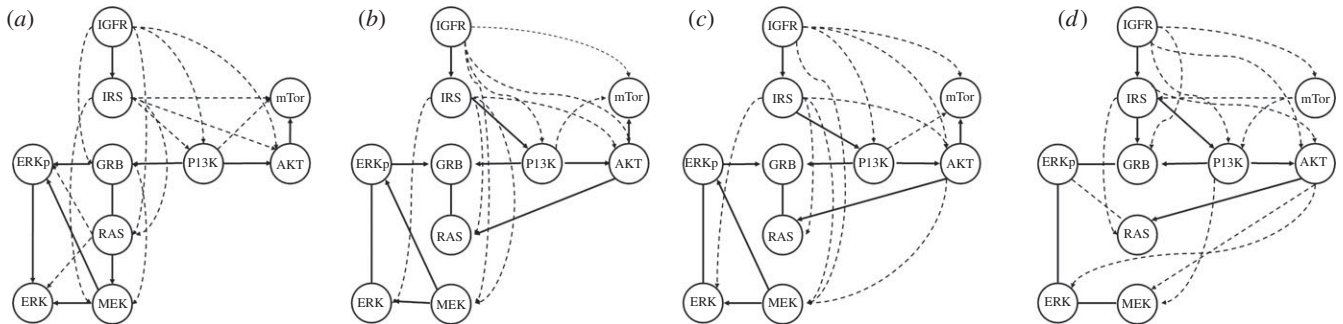


Figure 3. BN structure learning with perturbational data at different timepoints: (a) $t = 1$, (b) $t = 2$, (c) $t = 3$ and (d) $t = 4$. Solid lines indicate accurate edges, dotted lines indicate non-causal edges.

10 extra edges in each result. While these could result from the unmodelled cycles, the number is surprisingly high, and it is difficult to imagine that two feedback loops could result in such a large number of non-causal edges. Confusion over these unexplained edges led to the theoretical exploration in §4. The models obtained are reasonably robust over time, so the choice of an optimal timepoint does not appear crucial in this system, as long as perturbational data are used. This has important (and reassuring) implications since in real biological systems, the snapshot model is typically learned on a timepoint which is chosen either arbitrarily, or on the basis of the intuition of the biology expert. As before, there does not appear to be an indication of early versus late events. Thus, even with perturbational data, it is not possible to infer dynamics from the timecourse of snapshot models.

4. Impact of dynamics on learned snapshot models

As detailed in §2.3, a large number of factors confound most structure-learning efforts. Our idealized synthetic data eliminate the majority of these effects: hidden variables are not present, perturbations are completely effective and perfectly selective, data are abundant and not excessively noisy. The presence of cycles is the lone remaining potential confounder. Can the large number of non-causal edges all be attributed to cycles in the system? Or is there an additional factor impacting these results?

Despite the prevalence of snapshot models in biology, there is no theoretical study establishing the basic assumption

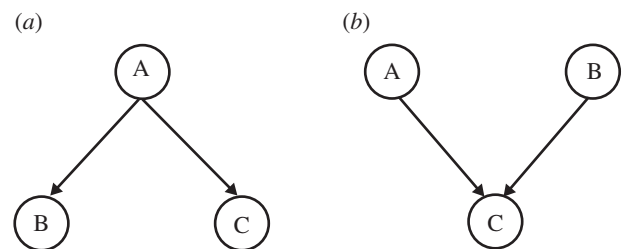


Figure 4. (a) True dynamic model and (b) statistical snapshot model of example 4.1, learned using observational data.

on which their success relies, i.e. if a dynamic system is measured at a single timepoint from cells or cell populations with randomly distributed initial conditions, then the joint distribution of the variables measured in this snapshot dataset can be used to elucidate causal relationships that exist among the variables and govern the behaviour of the dynamic system. We sought to prove the above statement, only to realize that it is not, in general, necessarily true. We illustrate this with an example.

Example 4.1

Consider three binary random processes A , B and C such that $A(0) = B(0) = C(0) = 0$, and for $t > 0$, $A(t)$ is sampled randomly and independently with probability 0.5 of being 0 or 1, $B(t) = A(t-1)$ and $C(t) = A(t) \vee A(t-1)$, where \vee is the binary addition or the OR function. The model is shown in figure 4.

The structure of the dynamic system has A as a root and B and C as its children. The snapshot picture (at any time $t > 1$), however, is very different since $B(t)$ is not even correlated with

$A(t)$, but $C(t)$ is determined by both of them. Thus, the structure of the conditional independencies is a V structure with A and B as the parents of C .

How does this pertain to dynamic systems in biology? Anytime a parent variable and its descendants are affected at different rates by the dynamics of the system, the dynamics can serve as a confounding effect. Non-causal edges result whenever an additional variable represents the abundance of a causal parent at the time of parent to child impact, better than the measured abundance of the parent variable itself at the time of measurement (i.e. at the time the snapshot is recorded). In the example above, since $B(t)$ has useful information about the history of $A(t)$ ($A(t-1)$ in this case) that $A(t)$ does not have, $B(t)$ can help predict $C(t)$ at any measured timepoint. This non-causal statistical dependency confounds structure-learning efforts, and probably represents the main culprit behind the abundance of non-causal edges that appear in our results above.

4.1. Impact of perturbations

Let us consider the effect of perturbations on this putative confounding effect. A perturbation on the ‘activity’ of $A(t)$ that makes it feed into B and C as 1 (independently of what $A(t)$ really is) will reflect the correct structure: The distribution of the snapshot data of B and C at any timepoint will differ from the unperturbed case. Moreover, a similar perturbation on B would not affect the snapshot data of A or C , which show that the V structure recovered by the statistical methods does not reflect the causal dynamic structure. Note that since BNs choose graphs that do well under all conditions used in learning, the confounding effects can be expected to dissipate given exhaustive or sufficient perturbations. However, with a reasonable number of perturbations applied above, our results nonetheless included a large number of non-causal edges.

Despite this observation, it would nonetheless seem that it should be possible to improve the quality of structure recovery using perturbations, since the effects of perturbations are usually more faithful to the dynamic system. To see this, we consider a dynamic system S with variables x_1, \dots, x_n , random initial conditions $X_1(0), \dots, X_n(0)$ and a set of perturbations I_1, \dots, I_m , and introduce a notion of *sensitivity* to our models:

Definition 4.2. *Sensitivity:* S is called sensitive at time $T > 0$ if $\forall t: 0 < t \leq T$, X_i and pairs of perturbations I_j and I_k , if

$$P_k[\pi(X_i)(t)] \neq P_j[\pi(X_i)(t)],$$

where P_j is the probability distribution under perturbation I_j , then

$$P_k[X_i(T)] \neq P_j[X_i(T)],$$

where $\pi(X_i)$ is the parent set of (X_i) in the graph structure. Sensitivity merely says that changes in the dynamic parents of a variable change that variable, or in other words, that each variable is sensitive to changes in its parents; when the parents change, the child node responds by changing as well. It can be seen that given a mass action model of a dynamic system, sensitivity holds true (for some range of T) almost surely. Sensitive systems satisfy the following proposition.

Proposition 4.3. *Given a system S that is sensitive at time T and has variables x_1, \dots, x_n with random initial conditions, and a set of perturbations I_1, \dots, I_m on x_1, \dots, x_n , then*

$$P_j[X_i(T)] \neq P[X_i(T)],$$

if and only if an ancestor of X_i was perturbed in I_j .

It is easy to show that the proposition holds between a variable and its parents; that is, given a system S that is sensitive at time T and has variables x_1, \dots, x_n with random initial conditions, and a set of perturbations I_1, \dots, I_m on x_1, \dots, x_n , then

$$P_j[X_i(T)] \neq P[X_i(T)],$$

if and only if the distribution of $\pi(X_i)$ was changed under I_j . From the definition of sensitivity, if

$$P_j[\pi(X_i)(t)] \neq P[\pi(X_i)(t)],$$

then

$$P_j[X_i(T)] \neq P[X_i(T)],$$

and therefore if the distribution of X_i 's parents change, then X_i 's distribution will change. Now assume that under perturbation I_j , the distribution of the parents of X_i did not change for all t such that $0 < t \leq T$. We can condition X_i on the history of its parents to see that its distribution does not change:

$$\begin{aligned} P_j(X_i(T)) &= P(X_i(T) | \pi(X_i(t))) P_j(\pi(X_i(t))) \\ &= P(X_i(T) | \pi(X_i(t))) P(\pi(X_i(t))), \end{aligned}$$

since

$$P_j(\pi(X_i(t))) = P(\pi(X_i(t))).$$

Now Proposition 4.3 holds by the iterative application of the above relationship: without loss of generality, we can re-enumerate the variables in S such that the perturbed variables are X_1, \dots, X_k , where $k < n$, and the variable of interest is X_n . By iterating the above relationship at most n times, we find that

$$P_j(D(X_1, \dots, X_k)) \neq P(D(X_1, \dots, X_k)),$$

where $D(X_1, \dots, X_k)$ is the set of all of the descendants of X_1, \dots, X_k . Thus if X_n is a descendant of any of X_1, \dots, X_k , then the distribution of X_n is modified.

Similarly, if none of X_n 's ancestors were perturbed under I_j , then the distribution of X_n does not change.

Based on the claims above, when a perturbation is employed, children of the target node will be (detectably) affected in a way that is consistent with the causal structure. This implies that reliance on the available perturbation data can ‘rescue’ the confounding effects of system dynamics, resulting in a more causal model structure, because any time a perturbation is used, it enforces correct causal edges. Between the causal and non-causal, confounded edges, there exists, it seems, a balance, with different edges outweighing each other depending on the relative strengths of the confounding versus causal effects in the system. It follows from the above that perturbational data may help tip the balance in the direction of causal edges. The perturbation data used in BN learning above was not sufficient to remove non-causal edges in the model, yet we hypothesized that it may be sufficient if used in a way that relies more heavily on the observed effects of perturbations. We therefore used the same data in our GBN model, presented below.

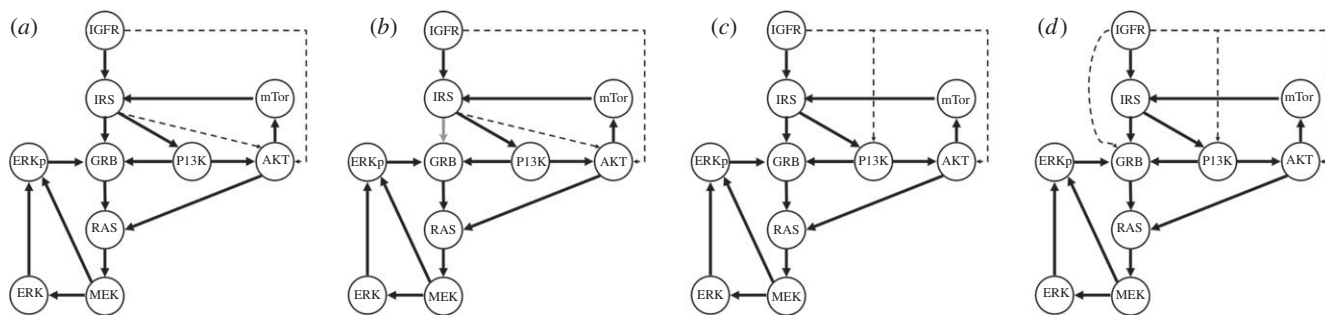


Figure 5. Graphs produced by the GBN algorithm at different timepoints: (a) $t = 1$, (b) $t = 2$, (c) $t = 3$ and (d) $t = 4$. Dotted edges are non-causal.

4.2. General Bayesian network learning

Based on the above study, perturbations improve the causal nature of learned models by minimizing the potentially confounding effects of system dynamics. Yet, standard BN learning of a 10 variable network with five unique perturbations yields approximately 10 non-causal edges, possibly due to persisting confounding effects. Feedback loops in the system may also confound the learned BN structure. In this section, we learn model structure with GBNs. We hypothesized that the GBNs would fare better than standard BNs for several reasons. First, they allow cycles, eliminating any connectivity errors induced by the acyclicity constraint. Next, they use perturbation information in a way that is qualitatively different from standard BNs: in GBN learning, descendants of perturbed variables are detected, and the resulting structure is constrained to contain a directed path between each discovered ancestor–descendent pair. Finally, structure learning takes place when all the cycles are broken, somewhat halting the dynamics of the system and potentially reducing any confounding effects. In figure 5, we present results for GBN structure learning, implemented as described in §2.2. Note that the data employed are identical to that used in standard BN learning.

The structure recovered is in general very close to the connectivity structure of the dynamic system and shows robustness over time, confirming the BN results showing lack of sensitivity to the specific timepoint selected. The improvement in elucidation of causal edges can be attributed to the ability to learn cycles, while the dramatic reduction in non-causal edges is probably due to the fact that descendent information is very robust (when the system is sensitive), allowing the GBN algorithm to avoid the confounding effects of the dynamics. There are some additional edges that seem to be persistent at all timepoints, two to three edges in each graph. The reason for those correlations is probably persistent confounding effects or some unrepresented effects in the system.

5. Conclusions

In this work, we have studied the relationship between the causal structure of a dynamic system and the reflected structure in snapshot data. We studied the situations in which single timepoint data can give misleading results, and the cases in which perturbational data can help in the recovery of the true structure, for reasons *other* than their well-documented utility in orienting edges. We tested our results on an accurate dynamic model of the IGF signalling pathway, and showed that it is possible to recover the structure of the dynamic causal relations from static data. We also show that when

models are learned at separate timepoints, the results are generally insensitive to the timepoint selected, and that the dynamics of the system are generally not revealed.

In our structure-learning efforts, we found that the correct connectivity structure is not fully or exclusively recovered by standard BN learning, even given a highly optimal dataset including many measurements, ideal perturbations (fully effective and selective), and a large number of perturbations. One important reason for this may be non-causal correlations imposed by the dynamics of the system. This phenomenon has previously been described for dynamic models [16–18], but it has, to our knowledge, never previously been considered in the context of static models. Therefore, this work brings to light an important yet hitherto unacknowledged confounder of structure inference efforts, and provides insight towards potential solutions. The understanding that underlying dynamics may confound structure elucidation for snapshot models can have far-reaching implications for this important goal in computational molecular biology.

The confounders presented here may be handled very differently if dynamic data were used, especially if fine time resolution data were available. The discussion on dynamic data are beyond the scope of this work, which focuses specifically on snapshot models, but we note that at least some aspects of this problem would not necessarily be eliminated, in part because some effects occur in infinitesimally small time slices. While the results presented are for a specific dynamic system, it is a system with components that appear frequently in biological networks, both at the node level and at the level of the types of feedback displayed. Therefore, we believe these results will reflect on many systems. Other limitations of structure learning, such as insufficient data, inadequate perturbations and the presence of hidden variables can also contribute to non-causal structure learning; whether these effects dominate over the one we have now described remains an open question. The answer is likely to depend on the particular system under study.

GBNs improve substantially upon the BN results, due in part to their fuller use of perturbation information. While the power of perturbations to direct edges and minimize the confounding effects of hidden variables is well established, this work emphasizes a *separate contribution* of perturbation information: elimination of confounding effects for children of the perturbed variable. This effect is related to hidden variables, in that non-causal edges in the structure-learning results are due to correlations of a variable with non-causal parents which carry information on earlier timepoints of the variable's causal parents, information which is not contained in the causal parent itself, at the current timepoint. In other words, the *history* of the variable's *non-hidden* parents is in

essence the hidden variable that induces non-causal edges. This is an unconventional form of hidden variables, which are normally assumed to be additional entities participating in the system (at the given timepoint); but their effects are alleviated similarly, with perturbations. Appropriate use of this information can enable elucidation of increasingly causal models of dynamic systems.

Endnote

¹It is difficult to assess what is a sufficient dataset size for structure learning. However, based on experience we can surmise that a dataset of 1000 is of sufficient size to learn a network of this size and

complexity. This is supported by the fact that high scoring models strongly outperform (outscore) lower scoring models (data not shown).

This work was supported by the NCI ICBP U54CA112970 and U54CA149145 grants. K.S. was supported by a Leukemia and Lymphoma Society Fellowship. G.P.N. is supported by the Rachford and Carlota A. Harris Endowed Professorship and grants from U19 AI057229, P01 CA034233, HHSN272200700038C, 1R01CA130826, CIRM DR1-01477 and RB2-01592, NCI RFA CA 09-011, NIH41000411217, NHLBIHVN01-HV-00242, European Commission HEALTH.2010.1.2-1, Merck, Novo Nordisk Biotech, Gilead Sciences, Inc., Celgene, Alliance for Lupus Research, DODCDMRP, NIH5-24927, FDA BAA-12-00118, and the Entertainment Industry Foundation NWCRA. C.J.T. is additionally supported by the NCI PSOC U54CA143826 grant.

References

- Demeure MJ, Bussey KJ, Kirschner LS. 2011 Targeted therapies for adrenocortical carcinoma: IGF and beyond. *Horm. Cancer* **2**, 385–392. (doi:10.1007/s12672-011-0090-6)
- Gately K, Collins I, Forde L, Al-Alao B, Young V, Gerg M, Feuerhake F, O'Byrne K. 2011 A role for IGF-1R-targeted therapies in small-cell lung cancer? *Clin. Lung Cancer* **12**, 38–42. (doi:10.3816/CLC.2011.n.005)
- Sachs K, Perez O, Pe'er D, Lauffenburger DA, Nolan GP. 2005 Causal protein-signaling networks derived from multiparameter single-cell data. *Science* **308**, 523–529. (doi:10.1126/science.1105809)
- Friedman N, Lital M, Nachman I, Pe'er D. 2000 Using Bayesian networks to analyze expression data. *J. Comput. Biol.* **7**, 3–4. (doi:10.1089/106652700750050961)
- Woolf PJ, Prudhomme W, Daheron L, Daley G, Lauffenburger DA. 2004 Bayesian analysis of signaling networks governing embryonic stem cell fate decisions. *Bioinformatics* **21**, 741–753. (doi:10.1093/bioinformatics/bti056)
- Dhaeseleer P, Wen X, Fuhrman S, Somogyi R. 1998 Mining the gene expression matrix: inferring gene relationships from large scale gene expression data. *IPCA* **98**, 203–212.
- Margolin A, Wang K, Califano A, Nemenman I. 2010 Multivariate dependence and genetic networks inference. *IET Syst. Biol.* **4**, 428. (doi:10.1049/iet-syb.2010.0009)
- Boscolo R, Liao JC, Roychowdhury VP. 2008 An information theoretic exploratory method for learning patterns of conditional gene coexpression from microarray data. *IEEE/ACM Trans. Comp. Biol. Bioinf.* **5**, 15–24.
- Butte AJ, Kohane IS. 2000 Mutual information relevance networks. *Pacific Symp. Biocomput.* **2000**, 418–429.
- Margolin AA, Nemenman I, Basso K, Wiggins G, Stolovitzky G, Favera RD, Califano A. 2006 ARACNE: an Algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinf.* **7**(Suppl 1), article S7.
- Friedman N, Murphy K, Russell S. 1998 Learning the structure of dynamic probabilistic networks. In *Proc. Fourteenth Annu. Conf. on Uncertainty in Artificial Intelligence, Madison, WI, 24–26 July 1998*, pp. 139–147. San Francisco, CA: Morgan Kaufmann.
- Sachs K, Itani S, Carlisle J, Nolan GP, Pe'er D, Lauffenburger DA. 2009 Learning signaling network structures from sparsely distributed data. *J. Comput. Biol.* **16**, 201–212. (doi:10.1089/cmb.2008.07TT)
- Sachs K, Gentles AJ, Youland R, Itani S, Irish J, Nolan GP, Plevritis SK. 2009 Characterization of patient specific signaling via augmentation of Bayesian networks with disease and patient state nodes. *IEEE Eng. Med. Biol. Soc.* **1**, 6624–6627.
- Dawid AP. 2008 Beware of the DAG!. In *JMLR: Workshop and Conference Proceedings* **6**, 59–86.
- Dash D. 2003 Caveats for causal reasoning with equilibrium models. PhD thesis, Intelligent Systems Program, University of Pittsburgh, USA.
- Kolar M, Song L, Ahmed A, Xing EP. 2010 Estimating time varying networks. *Ann. Appl. Statist.* **4**, 94–123. (doi:10.1214/09-AOAS308)
- Nodelman U, Shelton C, Koller D. 2002 Continuous time Bayesian networks. In *Proc. Eighteenth Conf. Uncertainty in Artificial Intelligence (UAI)*, pp. 378–387.
- Nodelman U, Shelton C, Koller D. 2003 Learning continuous time Bayesian networks. In *Proc. Nineteenth Conf. Uncertainty in Artificial Intelligence (UAI)*, pp. 451–458.
- Itani S, Ohannessian M, Sachs K, Nolan GP, Dahleh MA. 2008 Structure learning in causal cyclic networks. *J. Mach. Learn. Res. Workshop Conf. Proc.* **6**, 165–176.
- Sachs K, Itani S, Dahleh MA, Nolan GP. 2009 Learning cyclic signaling pathway structures while minimizing data requirements. *Pacific Symp. Biocomput.* **14**, 63–74.
- Pearl J. 1988 *Probabilistic reasoning in intelligent systems*. Las Altos, CA: Morgan Kaufmann.
- Hartemink AJ, Gifford DK, Jaakkola TS, Young RA. 2001 Using graphical models and genomic expression data to statistically validate models of genetic regulatory networks. *Pacific Symp. Biocomput.* **2001**, 422–433.
- Sachs K, Gifford D, Jakkola T, Sorger P, Lauffenburger DA. 2002 Bayesian network approach to cell signaling pathway modeling. *Sci. Signal.* **148**, PE38. (doi:10.1126/scisignal.1482002pe38)
- Pearl J, Verma TS. 1991 A theory of inferred causation. In *Principles of knowledge representation and reasoning: Proc. of the Second Int. Conf. (eds J Allen, R Fikes, E Sandewall)*, pp. 441–452. San Mateo, CA: Morgan Kaufmann.
- Heckerman D, Meek C, Cooper GF. 1999 A Bayesian approach to causal discovery. In *Computation, causation, and discovery (eds C Glymour, GF Cooper)*, pp. 141–166. Cambridge, MA: MIT Press.
- Bollen KA. 1989 *Structural equations with latent variables*. New York, NY: Wiley.
- Schmidt M, Murphy K. 2009 Modeling discrete interventional data using directed cyclic graphical models. In *Proc. 25th Annu. Conf. Uncertainty in Artificial Intelligence*, pp. 487–495.
- Carlson CJ. 2004 Mammalian target of rapamycin regulates IRS-1 serine 307 phosphorylation. *Biochem. Biophys. Res. Comm.* **316**, 533–539. (doi:10.1016/j.bbrc.2004.02.082)
- Moelling K, Schad K, Bosse M, Zimmermann S, Schwenecker M. 2002 Regulation of Raf-Akt Crosstalk. *J. Biol. Chem.* **277**, 31 099–31 106. (doi:10.1074/jbc.M111974200)