

Published in IET Systems Biology  
 Received on 1st September 2007  
 doi: 10.1049/iet-syb.2008.0130



# Genetic network identification using convex programming

A. Julius<sup>1</sup> M. Zavlanos<sup>1</sup> S. Boyd<sup>2</sup> G.J. Pappas<sup>1</sup>

<sup>1</sup>Department of Electrical and Systems Engineering, University of Pennsylvania, PA, USA

<sup>2</sup>Department of Electrical Engineering, Stanford University, CA, USA

E-mail: agung@seas.upenn.edu

**Abstract:** Gene regulatory networks capture interactions between genes and other cell substances, resulting in various models for the fundamental biological process of transcription and translation. The expression levels of the genes are typically measured as mRNA concentration in micro-array experiments. In a so-called genetic perturbation experiment, small perturbations are applied to equilibrium states and the resulting changes in expression activity are measured. One of the most important problems in systems biology is to use these data to identify the interaction pattern between genes in a regulatory network, especially in a large scale network. The authors develop a novel algorithm for identifying the smallest genetic network that explains genetic perturbation experimental data. By construction, our identification algorithm is able to incorporate and respect a priori knowledge known about the network structure. A priori biological knowledge is typically qualitative, encoding whether one gene affects another gene or not, or whether the effect is positive or negative. The method is based on a convex programming relaxation of the combinatorially hard problem of  $L_0$  minimisation. The authors apply the proposed method to the identification of a subnetwork of the SOS pathway in *Escherichia coli*, the segmentation polarity network in *Drosophila melanogaster*, and an artificial network for measuring the performance of the method.

## 1 Introduction

The use of RNA microarray has made it possible to have an expression profile for a large number of genes when exposed to different conditions. One of the most important problems in systems biology is to use these data to identify the interaction pattern between genes in a regulatory network, especially in a large scale network. In the literature, this is sometimes called reverse engineering the genetic network. Genetic network identification has important potential applications, for example, in drug discovery where a systems wide understanding of the regulatory network is crucial for identifying the targeted pathways.

Genetic network identification is a very active research field. For an overview on existing results and methodologies, we refer the reader to [1–5] and the references therein. On the basis of the technique, we identify two classes of methods for network identification. The first class consists of methods that infer

the network by clustering the genes based on their expression profiles. These methods do not view the network as a dynamical system, and the inferred network typically lacks causality.

The second class consists of methods that infer the cause-effect relation between genes through various representations. Several typical representations are information-theoretic network Bayesian network and dynamical network described by ordinary differential equations (see the survey in [5]). Information-theoretic network-based methods typically lack the causality information, as they identify the network as undirected graphs. Bayesian network-based methods identify the genetic network as directed graph and thus convey some causality information. However, they typically do not accommodate cycles in the network graph. This limitation can be significant, as feedback motifs are very common in genetic regulatory networks. Both causality and feedback motives limitation are not present in methods that model the

network as the network is modelled as a set of differential equations [6–9]. The method that we propose in this paper belongs to this class.

On the basis of the type of data used in the identification, there are two classes of methods. The first class deals with data obtained from dynamic time-series measurement of the expression profiles. The second class deals with steady-state data, obtained by measuring the expression profiles when the network reaches an equilibrium. Our method belongs to the second class, where the identification of the interconnection pattern is done locally by perturbing the network around a given equilibrium. It is generally known that a regulatory network can have multiple stable equilibria.

The method that we propose aims at providing a minimal model that explains given genetic perturbation data. Obtaining such a minimal model is computationally very hard, as it involves combinatorial exploration of all possible network topologies. Such a problem has been shown to have NP-Hard complexity [10]. Pioneering works by Collins and coworkers [6, 7, 11] provided an important step towards addressing this problem. In [11], the authors propose a method for noiseless measurements, where the optimisation is relaxed as a non-recursive  $\ell_1$  optimisation problem. In [6, 7], the method that they use introduces an a priori limitation on the connectivity of the network, and perform a combinatorial search on this limited set. The connectivity limitation is that each gene in the network has the same number of inputs. Another different approach where the connectivity is imposed on the number of outputs is reported in [12]. For every combination, the parameters of the model are deduced through a least-square fitting. Similar approach that uses least-square fitting but without minimisation of the model is also reported in [13].

In solving this problem, we take a different approach. Instead of imposing a connectivity limitation on the network, we do not have any limit on how many inputs each gene should have. The hard combinatorial problem is solved using a mathematical technique called convex optimisation [14] after relaxing it as a recursive  $\ell_1$  optimisation problem [15]. The same technique has been applied successfully in various fields where sparsity optimisation is needed such as, portfolio optimisation in finance [16] and controller design in engineering [17]. Different techniques of convex  $\ell_1$  relaxation have been used in other works in reverse engineering of gene networks, such as [18–22]. Posing the problem as a convex optimisation problem is actually very advantageous from the complexity point of view, as convex optimisation algorithms can be implemented reliably to solve large scale problems.

Solving the problem with convex  $\ell_1$  relaxation has an advantage of being able to handle noisy data, incorporate a priori knowledge about the network structure, encoding whether one gene affects another gene or not, or whether

the effect is positive or negative. The identified model is then constructed to satisfy the a priori knowledge by default.

In this paper, we apply our method to two networks that have been previously identified in the literature: the SOS pathway in *Escherichia coli* [16] and the segmentation polarity network in *Drosophila melanogaster* [7]. We also apply our method to an artificial network to assess its performance, primarily in relation to other convex  $\ell_1$  relaxation methods.

## 2 Gene network modelling and identification

### 2.1 Notations

In this paper, we use the following matrix notation. If  $X$  is a matrix with  $n$  rows and  $m$  columns, we write  $X \in \mathbb{R}^{n \times m}$ . The symbol  $X_{ij}$  refers to the entry of  $X$  at the  $i$ th row,  $j$ th column. A single index such as  $X_j$  refers to the column vector corresponding to the  $j$ th column of  $X$ . The operator  $\mathbb{E}[X_j]$  is the probabilistic expectation of  $X_j$ . The operator  $\text{Var}[X_j]$  is the covariance of  $X_j$ .

A genetic regulatory network consisting of  $n$  genes in a genetic perturbation experiment can be modelled as a dynamical system [6, 7]. In general, such a model assumes the following form

$$\frac{d\hat{x}}{dt} = F(\hat{x}, \hat{y}, \hat{u}), \quad \hat{x} \in \mathbb{R}^n, \quad \hat{y} \in \mathbb{R}^n, \quad \hat{u} \in \mathbb{R}^p \quad (1)$$

$$\frac{d\hat{y}}{dt} = G(\hat{x}, \hat{y}) \quad (2)$$

where  $\hat{x}_i \in \mathbb{R}$  denotes the transcription activity (typically measured as mRNA transcript concentration) of gene  $i$  in the network,  $\hat{y}_i$  denotes the protein concentration of protein  $i$  and  $\hat{u}_i$  is the so called transcription perturbation. In very large networks, we can typically assume that not all genes can be perturbed in the experiment, resulting in  $p < n$ .

The functions  $F$  and  $G$  summarise the dynamics of transcription and translation, as well as factors such as the degradation and dilution of transcripts and proteins. Such nonlinear genetic networks can have multiple stable equilibria. Each equilibrium typically corresponds to a phenotypical state of the system. The dynamics close to a given equilibrium  $(\mathbf{x}_{\text{eq}}, \mathbf{y}_{\text{eq}})$  can be approximated by the set of linear differential equations

$$\frac{d\mathbf{x}}{dt} = \mathfrak{U}_{11}\mathbf{x} + \mathfrak{U}_{12}\mathbf{y} + \mathbf{u} \quad (3)$$

$$\frac{d\mathbf{y}}{dt} = \mathfrak{U}_{21}\mathbf{x} + \mathfrak{U}_{22}\mathbf{y} \quad (4)$$

where  $\mathbf{x} := \hat{\mathbf{x}} - \mathbf{x}_{\text{eq}}$  and  $\mathbf{y} := \hat{\mathbf{y}} - \mathbf{y}_{\text{eq}}$  [8, 13]. The matrices  $\mathfrak{U}_{ij}$ ,  $i, j = 1, 2$  are the linearisation of the dynamics near the equilibrium, while the vector  $\mathbf{u} \in \mathbb{R}^n$  represents the

effect of the perturbation inputs in the linear model. In the case that not all genes can be perturbed,  $\mathbf{u}$  is constrained in a subspace of  $\mathbb{R}^n$ . Given that the system is stable around the equilibrium  $(\mathbf{x}, \mathbf{y}) = (0, 0)$ , if  $\mathbf{u}$  is small enough, the system will move to a new equilibrium  $(\mathbf{x}, \mathbf{y})$ , for which

$$\mathcal{A}_{11}\mathbf{x} + \mathcal{A}_{12}\mathbf{y} + \mathbf{u} = 0 \quad (5)$$

$$\mathcal{A}_{21}\mathbf{x} + \mathcal{A}_{22}\mathbf{y} = 0 \quad (6)$$

We formulate a theory for the case when only the transcript concentrations are measured. In this case, we can flatten the transcription and translation layer of the regulatory system into an effective gene–gene regulatory network, by eliminating the protein concentration  $\mathbf{y}$ . We therefore obtain

$$\mathbf{y} = -\mathcal{A}_{22}^{-1}\mathcal{A}_{21}\mathbf{x} \quad (7)$$

$$(\mathcal{A}_{11} - \mathcal{A}_{12}\mathcal{A}_{22}^{-1}\mathcal{A}_{21})\mathbf{x} + \mathbf{u} = 0 \quad (8)$$

We then define the matrix

$$\mathbf{A} := (\mathcal{A}_{11} - \mathcal{A}_{12}\mathcal{A}_{22}^{-1}\mathcal{A}_{21}) \quad (9)$$

The matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  encodes the effective pairwise interactions between the individual genes in the network. In this representation, we have the relation

$$\mathbf{A}\mathbf{x} + \mathbf{u} = 0 \quad (10)$$

Let  $\mathbf{U} = [U_1 \dots U_m] \in \mathbb{R}^{p \times m}$  denote the stack matrix of the transcription perturbations for different  $m$  experiments and  $\mathbf{X} = [X_1 \dots X_m] \in \mathbb{R}^{n \times m}$  denote the stack matrix of the corresponding steady-state mRNA concentrations. In large networks or in cases where experiments are costly, we can typically assume that the experimental data set is smaller than the network size [2]. That is, we assume  $m < n$ .

By collecting all  $m$  experiments at steady state, the equilibrium conditions (10) can be written as

$$\mathbf{A}\mathbf{X} + \mathbf{U} = 0 \quad (11)$$

In (11), the matrices  $\mathbf{X}$  and  $\mathbf{U}$  are known, since they are measured (possibly with noise). The goal of the method we propose in this paper is to find unknown matrix  $\mathbf{A}$ , which models genetic network interactions and best explains the genetic perturbation experiments.

We aim at constructing a model that not only explains the perturbation data, but also incorporates some a priori knowledge about the system. Here a priori biological knowledge is typically qualitative, encoding whether one gene affects another gene or not, or whether the effect is positive or negative. This knowledge is then manifested as pre-specified signs of some entries of the matrix  $\mathbf{A}$ . Furthermore, we would like to develop a model that constitutes a minimal network.

We quantify the size of the model as the number of connections in the network, that is, the number of non-zeros entries in the matrix  $\mathbf{A}$ . This is known as the  $L_0$  norm of the matrix  $\mathbf{A}$ . (Even though it is not a norm in a strict mathematical sense, it is common to refer to it as the  $L_0$  norm.) The minimal model then corresponds to a network with as few connections as possible, or equivalently, the sparsest possible matrix  $\mathbf{A}$ .

Obtaining a minimal model is clearly beneficial, as it reduces the complexity the model. A non-minimal model might be able to explain the data slightly better than a minimal one, for example, due to measurement noise. However, this can lead to a phenomenon called overfitting, where a model includes unnecessary features to accommodate the noise. A minimal model is also desirable when the identified model is used in a pathway knockout. A non-minimal model might contain falsely identified spurious pathways.

### 3 Identification algorithm

Consider (11), if we assume that the measurements are noise free and if we have a sufficient number  $m = n$  of independent experiments,  $\mathbf{X}$  is an invertible matrix, and we could obtain  $\mathbf{A}$  using  $\mathbf{A} = -\mathbf{U}\mathbf{X}^{-1}$ . However, as genetic networks grow dramatically in size as we reach genome-scale networks, having  $n$  independent experiments can be costly, both financially and timewise. Furthermore, the absence of noise is not a realistic assumption, except when we are dealing with *in silico* model. Consequently, we are not going to assume that right hand side of (11) is zero. Instead, we define identification error as

$$\boldsymbol{\eta} := \mathbf{A}\mathbf{X} + \mathbf{U} \quad (12)$$

and try to minimise  $\boldsymbol{\eta}$  (as a function of  $\mathbf{A}$ ) with respect to some metric, while obtaining a minimal model for  $\mathbf{A}$  and satisfying the a priori constraints that might be imposed on  $\mathbf{A}$ .

#### 3.1 Error criterion

In this paper, we use the total squared error as the error criterion

$$\text{Err} = \sum_{j=1, \dots, m} \sum_{i=1, \dots, n} \eta_{ij}^2 \quad (13)$$

However, if the covariance of the error in the  $j$ th experiment is known, then the sum can be replaced by a more accurate weighted sum

$$\text{Err} = \sum_{j=1, \dots, m} \sum_{i=1, \dots, n} \sum_{k=1, \dots, n} \eta_{ij} \eta_{kj} \mathbf{R}_{ik}^j = \sum_{j=1, \dots, m} \boldsymbol{\eta}_j^T \mathbf{R}^j \boldsymbol{\eta}_j \quad (14)$$

where  $\mathbf{R}^j$  is the inverse of the error covariance matrix in the  $j$ th experiment. The intuition behind this weight is that the identification error in the experiments with more

reliable data (smaller variance) weights more than that coming from less reliable data.

### 3.2 Model minimality criterion

We define the size of a model given by the matrix  $A$  as the number of non-zero entries in  $A$ , which is denoted by  $\|A\|_0$ . This is the number of connections in the model.

### 3.3 Priori knowledge constraint

Such knowledge typically has the form of (partial) sign pattern of the matrix  $A$ . We encode this by a matrix  $S \in \{0, +, -, ?\}^{n \times n}$ , where

$$\begin{bmatrix} S_{ij} = + \\ S_{ij} = - \\ S_{ij} = 0 \\ S_{ij} = ? \end{bmatrix} \Leftrightarrow \begin{bmatrix} A_{ij} \geq \varepsilon \\ A_{ij} \leq -\varepsilon \\ -\frac{\varepsilon}{2} \leq A_{ij} \leq \frac{\varepsilon}{2} \\ A_{ij} \in \mathbb{R} \end{bmatrix} \quad (15)$$

Here  $\varepsilon$  is a small number, below which a connection is considered negligible. In this paper, we set  $\varepsilon = 10^{-3}$ . In short, such a pattern encodes known positive interactions (+), negative interactions (-), the absence of interactions (0), or simply lack of knowledge (?) between any two genes in the network. For example, a matrix consisting of only (?) indicates no a priori knowledge about the network. Hereafter, we shall denote the set of all matrices  $A$  that satisfy a given a priori knowledge constraint  $S$  as  $\mathcal{S}$ . A critical property of the set  $\mathcal{S}$  that can be easily shown is that it is convex [14].

### 3.4 Bias due to mRNA decay

As defined in (9), the effective gene–gene network model  $A$  is comprised of two parts. The direct transcript–transcript factor  $\mathfrak{U}_{11}$  and the regulation through protein factor  $-\mathfrak{U}_{12}\mathfrak{U}_{22}^{-1}\mathfrak{U}_{21}$ .

We can typically assume that the dynamics of the concentration of protein  $i$ ,  $y_i$ , is determined solely by its own transcript concentration  $x_i$  and its decay and/or dilution process. In this case, both  $\mathfrak{U}_{21}$  and  $\mathfrak{U}_{22}$  in (7) are diagonal matrices with positive and negative entries, respectively. We can therefore write

$$-\mathfrak{U}_{22}^{-1}\mathfrak{U}_{21} =: A_1 \quad (16)$$

where  $A_1$  is a diagonal matrix with positive entries. Similarly, if we assume that the dynamics of the transcript concentration,  $x_i$ , is due to decay, we can write

$$\mathfrak{U}_{11} =: -A_2 \quad (17)$$

where  $A_2$  is another diagonal matrix with positive entries. We then have

$$A = A_1\mathfrak{U}_{12} - A_2 \quad (18)$$

where  $\mathfrak{U}_{12}$  represents transcription regulation by proteins. Consequently, any network model that results from identification using genetic perturbation data will have a negative bias on the diagonal terms. When the decay rates  $A_2$  are known, which is the case for the *in silico* model of the *Drosophila* segmentation polarity network, we can remove the bias. Otherwise, we have to take this into account when making statements about autoregulation of genes.

### 3.5 Convex programming

The method that we propose in this paper is based on a mathematical technique called convex programming. Basically, convex programming is a mathematical theory for minimisation of a convex cost function over a convex set of feasible solutions. Formulating the identification problem as convex programming is attractive because there are techniques for solving convex programming problems efficiently; see, example, [14].

### 3.6 Convex optimisation solver

The convex optimisation problems that we pose in this paper are solved using MATLAB with the toolbox *cvx* [23] running on an Intel Xeon 2.8 Ghz processor with 4 GB RAM. Here *cvx* makes forming and solving the problem easy, but at the cost of efficiency. However, custom-made implementations of convex optimisation algorithms can easily handle problems with thousands of variables allowing us in the future to handle genome scale problems.

The method that we use in this paper can be explained in two steps.

**3.6.1 Step 1 – establishing baseline error level:** In this step, we establish the least error level that a model can attain, while disregarding the model minimality criterion. As discussed above, when we assume no a priori knowledge about the statistics of the error, we can simply use the total squared error as the error criterion. Consequently, finding the baseline error level  $E_{bs}$  amounts to solving the following convex programming problem

$$\begin{aligned} &\text{minimise} \quad \sum_{j=1, \dots, m} \sum_{i=1, \dots, n} \eta_{ij}^2 \\ &\text{subject to} \quad \eta = AX + U, A \in \mathcal{S} \end{aligned} \quad (19)$$

with optimisation variable  $A$ .

When the statistics of the measurement errors in each experiment for  $X$  and  $U$  are known, we can compute the associated covariance of the error criterion as

$$\text{Var}[\eta_j] = A\text{Var}[X_j]A^T + \text{Var}[U_j] \quad (20)$$

As discussed above, the error criterion that we use is given in (14), where

$$R^j = (\text{Var}[\eta_j])^{-1} \quad (21)$$



Consequently, finding the baseline error level  $E_{bs}$  amounts to solving the following optimisation problem

$$\begin{aligned} & \text{minimise} \quad \sum_{j=1, \dots, m} \eta_j^T R^j \eta_j \\ & \text{subject to} \quad \eta = \mathbf{A}\mathbf{X} + \mathbf{U}, \mathbf{A} \in \mathcal{S}, \\ & \quad \mathbf{R}^j = (\mathbf{A}\text{Var}[\mathbf{X}_j]\mathbf{A}^T + \text{Var}[\mathbf{U}_j])^{-1} \end{aligned} \quad (22)$$

where  $\mathbf{A}$  is the variable.

However, this formulation is not convex. In order to solve it efficiently, we relax the problem by approximating the covariance matrices. First, assuming that the covariance matrices are identity matrices, we find the best model that minimises the error criterion (19). Denote this model as  $\tilde{\mathbf{A}}$ . The weight matrices  $\mathbf{R}^j$  are then given by

$$\mathbf{R}^j = (\tilde{\mathbf{A}}\text{Var}[\mathbf{X}_j]\tilde{\mathbf{A}}^T + \text{Var}[\mathbf{U}_j])^{-1} \quad (23)$$

The baseline error level  $E_{bs}$  is then computed by solving the following convex optimisation problem.

$$\begin{aligned} & \text{minimise} \quad \sum_{j=1, \dots, m} \eta_j^T \mathbf{R}^j \eta_j \\ & \text{subject to} \quad \eta = \mathbf{A}\mathbf{X} + \mathbf{U}, \mathbf{A} \in \mathcal{S} \end{aligned} \quad (24)$$

with  $\mathbf{A}$  as the variable.

**3.6.2 Step 2 – minimising the model:** The baseline error level and the approximated error covariance matrix that we obtain in Step 1 above are used in finding a minimal model that can explain the data reasonably well. That is, we search for a minimal model that results in an error level of at most  $\beta E_{bs}$ , where  $\beta \geq 1$  is a predetermined parameter. The bigger the  $\beta$ , the more variation we allow for the identified model, and thereby possibly obtain a smaller model at the cost of higher error level. Thus,  $\beta$  allows us to control the trade-off between model accuracy and model minimality.

Mathematically, Step 2 can be formulated as the following optimisation problem

$$\begin{aligned} & \text{minimise} \quad \|\mathbf{A}\|_0 \\ & \text{subject to} \quad \eta = \mathbf{A}\mathbf{X} + \mathbf{U}, \mathbf{A} \in \mathcal{S} \\ & \quad \sum_{j=1, \dots, m} \eta_j^T \mathbf{R}^j \eta_j \leq \beta E_{bs} \end{aligned} \quad (25)$$

where  $\mathbf{A}$  is the variable. We denote the solution of this problem as  $\mathbf{A}_{min}$ . Although the constraints in the problem above defines a convex feasible set, the cost function itself is not convex. In fact, the problem has combinatorial complexity as we have to search in the set of all possible interconnection patterns. This means that the complexity of the problem increases very rapidly with its size, and thus makes it practically impossible to solve it on a large scale. In order to solve this problem with convex programming, we relax the  $L_0$  minimisation problem as a recursive weighted  $\ell_1$  minimisation.

*Step 2.1:* Initiate  $\mathbf{A}_{old} = 0$  and  $W_{ij} = 1, i = 1, 2, \dots, n, j = 1, 2, \dots, n$

*Step 2.2:* Find  $\mathbf{A}_{update}$  by solving the convex optimisation problem

$$\begin{aligned} & \text{minimise} \quad \sum_{i=1, \dots, n} \sum_{j=1, \dots, n} W_{ij} |A_{ij}| \\ & \text{subject to} \quad \eta = \mathbf{A}\mathbf{X} + \mathbf{U}, \mathbf{A} \in \mathcal{S} \\ & \quad \sum_{j=1, \dots, m} \eta_j^T \mathbf{R}^j \eta_j \leq \beta E_{bs} \end{aligned} \quad (26)$$

where  $\mathbf{A}$  is the variable.

*Step 2.3:* Update  $W_{ij} = f(\mathbf{A}_{update, ij})$ . Here  $f(\cdot)$  is a function that assigns a weight matrix for the convex cost function in Step 2.2. The choice for the function is explained in more detail later.

*Step 2.4:* If  $\|\mathbf{A}_{old} - \mathbf{A}_{update}\|_F \geq \varepsilon$ , then update  $\mathbf{A}_{old} = \mathbf{A}_{update}$  and go to step 2.2, otherwise stop the iteration and return the current solution as the optimal value

$$\mathbf{A}_{min} = \mathbf{A}_{update} \quad (27)$$

Here  $\varepsilon$  is a small number that we choose to indicate the convergence of the iteration. Throughout this paper, we use  $\varepsilon = 10^{-3}$ .

## 3.7 Choosing the weight function

The weight function  $f(A_{ij})$  is designed such that the entries of  $\mathbf{A}$  that are small are given larger weight than the larger entries [15]. This is because we are interested in maximising the number of zero entries in  $\mathbf{A}$ . We thus emphasise more on reducing the entries that are already small. The generic form of  $f(A_{ij})$  that we use in this paper is as follows

$$f(A_{ij}) = \frac{\delta^p}{\delta^p + |A_{ij}|^p} \quad (28)$$

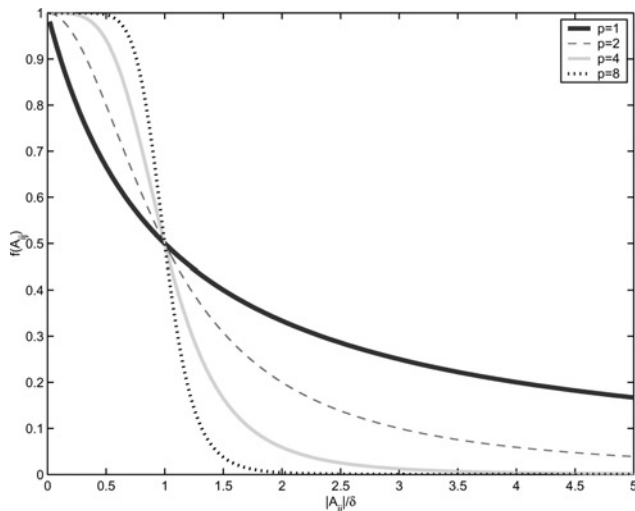
where  $\delta$  is a small number that acts as a threshold, below which a number is considered ‘small’. The parameter  $\delta$  thus determines the ‘boundary’ between small values and large values. The term ‘boundary’ is to be interpreted loosely, as the transition is smooth. The exponent  $p$  determines the shape of the function, and how abruptly the weight transitions take place from 1 to 0. The plot of the function  $f(A_{ij})$  can be seen in Fig. 1. Throughout this paper, we use  $\delta = 10^{-2}$  and  $p = 1$ .

The overall identification procedure can be summarised in the flowchart in Fig. 2.

## 4 Results and discussion

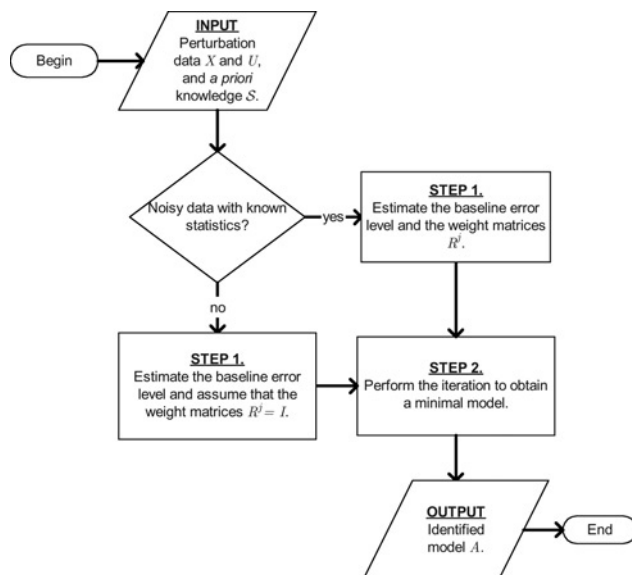
### 4.1 The segmentation polarity network of *Drosophila melanogaster*

We apply our proposed method to genetic perturbation data obtained from an in numero experiment based on the model



**Figure 1** The plots of the weight function  $f(A_{ij})$  against  $|A_{ij}|/\delta$  for various shape parameters  $p$

Notice that scaling the factor  $\delta$  amounts to shifting the transition between high weight and the low weight

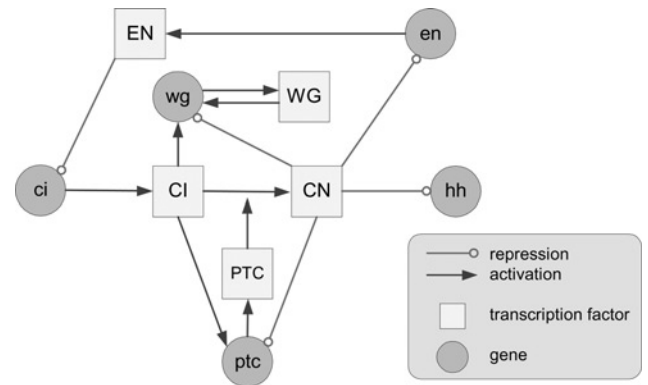


**Figure 2** The flowchart summarising the method proposed in this paper

The inputs to the algorithm are the perturbation data represented by the matrices  $X$  and  $U$ , as well as potential a priori structural knowledge about the network, which is represented by  $S$ . The output of the algorithm is the identified model, represented by the matrix  $A$ .

given in [7]. The network model captures the interaction between five genes and five transcription factors (Fig. 3).

The gene *ci* produces the *Cubitus interruptus* protein that further undergoes a post-translational modification into the activator form (CI) or the repressor form (CN). *Cubitus interruptus* activator acts as an activator for the genes *ptc* and *wg*, while the *Cubitus interruptus* repressor represses the genes *ptc*, *wg*, *en* and *hh*. The gene *wg* produces the



**Figure 3** The segmentation polarity network of *Drosophila melanogaster* [7]

protein wingless (WG) that in turn acts as a self-activator. The gene *ptc* produces the protein patched (PTC) that inhibits the modification of *Cubitus interruptus* into the activator form and promotes the modification into the repressor form, effectively forming a negative feedback loop. The gene *en* produces the protein engrailed (EN) that represses the transcription of *ci* and promotes the transcription of *hh*.

We obtain the perturbation data  $X$  and  $U$  by numerically integrating the model provided in [7]. We first simulate the model without any perturbation to obtain an equilibrium. We then perform simulations with constant perturbation to each of the five genes in the model. The perturbation is set at  $10^{-3}$ . Thus,  $U = 10^{-3} \cdot I$ , where  $I$  is an identity matrix. The deviations from the unperturbed equilibrium are recorded in the  $X$  matrix.

Since the data is noiseless, we obtain the baseline error level by solving (19). We do not estimate any error covariance matrices, and the performance measure is thus given by (13). We then proceed with Step 2, and use  $\beta = 1.1$ .

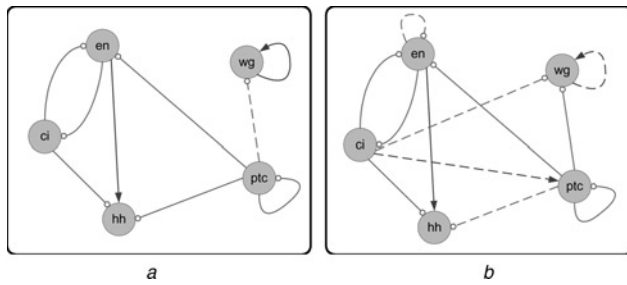
From the numerical model, we have precise knowledge about the mRNA decay rate. Using this information, we eliminate the negative (auto repression) bias in the identified network by adding a diagonal matrix  $A_2$  (18).

Given the matrix  $A$  of the identified model, we denote the strongest intergene interaction as

$$\mu := \max_{i \neq j} \|A_{ij}\| \quad (29)$$

An interaction represented by  $A_{ij}$  is considered strong if  $\|A_{ij}\| > 0.1 \cdot \mu$ , and weak if  $10^{-3} \leq \|A_{ij}\| \leq 0.1 \cdot \mu$ .

We execute our method to obtain a network model. The execution takes about 6 s on the platform that is detailed in the previous section. Fig. 4 shows the result of our network identification method and the network identified in [7]. Our model uses the full set of numerical data from the model, in



**Figure 4** Identified network using full numerical data from the model of the segmentation polarity network of *Drosophila melanogaster*

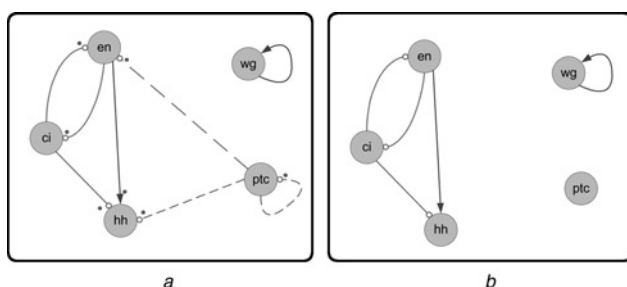
*a* The network identified using the method proposed in this paper

*b* The network identified in [7]

Solid lines indicate strong interaction, while broken lines indicate weak interaction

which all five genes are perturbed separately. We can see that our method produces a smaller model, and that all the identified connections can be accounted for based on the description of the network earlier in this section. The network in [7] contains a self-repression loop in *en*, which is a false positive as it is not included in the real model. The connections from *ci* to *wg* and *ptc* are not present in our model. As explained above, *ci* produces both an activator and repressor for *wg* and *ptc*. These connections are thus very weak and not included in our model.

We also apply the method on a partial data set. In this set, we do not include the data from the perturbation of *ptc*. The network identified using this data set is shown in Fig. 5*b*. Observe that the connections from *ptc* to other genes disappear as expected, since there is no data that dictates their existence. This lack of data can be supplemented by a priori knowledge. We then include a sign pattern matrix as a constraint. The result is shown in Fig. 5*a*. This network includes weak connections from *ptc* to itself, *en*, and *hh* as



**Figure 5** Identified network using partial numerical data from the model of the segmentation polarity network of *Drosophila melanogaster*

*a* The network identified when a priori knowledge about the sign pattern is incorporated

Arrows with black dots indicate interactions that are included in the a priori knowledge

*b* The network identified without incorporating any a priori knowledge

required by the constraint. These connections are weak because their existence is required by the a priori knowledge without any supporting data.

## 4.2 SOS pathway in *Escherichia coli*

We also apply our proposed method on a subnetwork of the SOS pathway in *Escherichia coli*, using the genetic perturbation experimental data set provided in [6]. The subnetwork that we consider consists of nine genes and several transcription factors and metabolites (see Fig. 2 in the Supplementary Material).

The main pathway featured in this network is the pathway between the single-stranded DNA (ssDNA) and the protein LexA that acts as a repressor to several other genes (*recA*, *ssb*, *dinI*, *umuDC* and *rpoD*). The protein RecA, which is activated by the ssDNA, cleaves LexA and thus upregulates the above mentioned genes. Other key regulators in the network are the sigma factors  $\sigma 70$ ,  $\sigma 32$  and  $\sigma 38$ . These sigma factors play an important role in initiating transcription in heat shock and starvation responses.

We obtain the perturbation data  $X$  from [6]. Since there is no explicit mentioning of  $U$ , we assume that it is an identity matrix. Notice that this is a justifiable assumption, as a different value of  $U$  would just result in a scaling of the model.

The covariance matrix of the measurement in  $X$  is obtained by processing the standard error matrix provided in [6]. Assuming that the measurement errors of different genes are uncorrelated, we can obtain

$$(\text{Var}[X_j])_{ik} = \begin{cases} \sigma_{ij}^2, & i = k \\ 0, & i \neq k \end{cases} \quad (30)$$

where  $\sigma_{ij}$  is the standard error of the measurement of gene  $i$  in the  $j$ th experiment. Since there is no information about the measurement error for  $U$ , we assume it is zero.

We compile the connections that are included in the a priori knowledge in Table 1. This list is compiled based on the diagram in Fig. 6. We begin with Step 1 of the method to obtain a baseline error level and estimated error covariance. We use the estimated error covariance to obtain the weight matrices  $R^j$  that are used in the error criterion (14). We then proceed to Step 2, and use different  $\beta$  values obtain different models and analyse them (as detailed below).

As a comparison, we also perform the identification without estimating the error covariance and using the weighted sum (14) as our error criterion. Instead, we use identity weight matrices. This approach turns out to be inferior to the one with estimated error covariance.

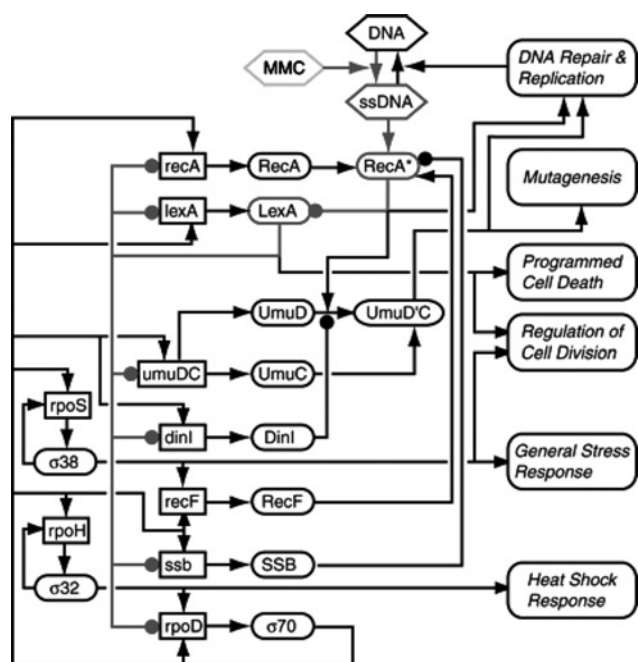
The result of our method is shown in Fig. 7. In panel (a), we see the network identified using our method with estimated error covariance and  $\beta = 1.75$ . The network

**Table 1** A summary of a priori knowledge used in the network identification

Genes	<i>recA</i>	<i>lexA</i>	<i>ssb</i>	<i>recF</i>	<i>dinI</i>	<i>umuDC</i>	<i>rpoD</i>	<i>rpoH</i>	<i>rpoS</i>
<i>recA</i>	?	—	?(-)	?(+)	?(+)	?(-)	+	?(0)	?(0)
<i>lexA</i>	+	—	?(-)	?(+)	?(+)	?(-)	+	?(0)	?(0)
<i>ssb</i>	+	—	?(-)	?(+)	?(+)	?(-)	+	?(0)	?(0)
<i>recF</i>	?(0)	?(0)	?(0)	?(-)	?(0)	?(0)	+	?(0)	+
<i>dinI</i>	+	—	?(-)	?(+)	?	?(-)	+	?(0)	?(0)
<i>umuDC</i>	+	—	?(-)	?(+)	?(+)	?(-)	+	?(0)	?(0)
<i>rpoD</i>	+	—	?(-)	?(+)	?(+)	?(-)	?	+	?(0)
<i>rpoH</i>	?(0)	?(0)	?(0)	?(0)	?(0)	?(0)	+	?	?(0)
<i>rpoS</i>	?(0)	?(0)	?(0)	?(0)	?(0)	?(0)	+	?(0)	?

The values in brackets represent known connections based on [6]. A + sign indicates known activation, — indicates known inhibition, 0 indicates the absence of connection and ? indicates unknown connection

identified in [6] is shown in Fig. 7 panel (b). Comparing it to our result in panel (a), we can see that the network in (b) misidentifies several known one-hop interconnections, such as the mutual repression between *recA* and *rpoD*.

**Figure 6** (Taken from [6]) The diagram of interactions in the SOS network

DNA lesions caused by mitomycin C (MMC) (light grey hexagon) are converted to single-stranded DNA during chromosomal replication. Upon binding to ssDNA, the RecA protein is activated (RecA\*) and serves as a coprotease for the LexA protein. The LexA protein is cleaved, thereby diminishing the repression of genes that mediate multiple protective responses. Boxes denote genes, ellipses and denote proteins, hexagons indicate metabolites, arrows denote positive regulation, filled circles denote negative regulation. Dark grey emphasis and denotes the primary pathway by which the network is activated after DNA damage

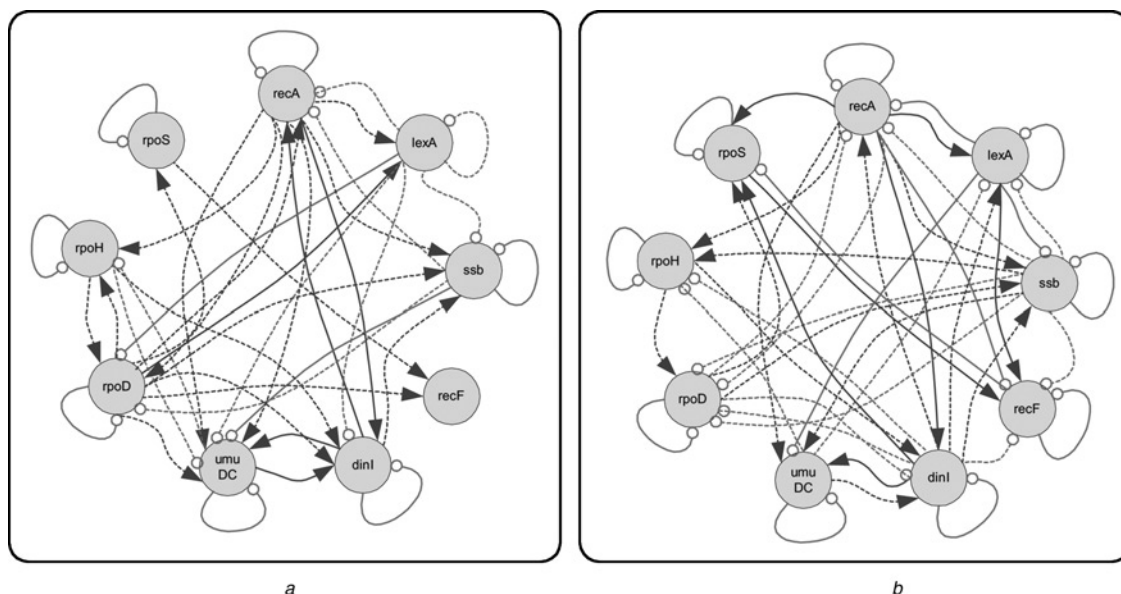
The network that we identify in panel (a) includes a number of interactions that are not included in the a priori knowledge. Upon cross-validation with the literature about the SOS network, we found that some of these new interactions are valid. For example, the protein *dinI* is known to stabilise *recA*\*, the activated form of *recA* [24]. Thereby, it effectively promotes the degradation of LexA, and thus activates *recA*, *lexA*, *ssb*, *dinI*, *umuDC* and *rpoD*. Our result correctly predicts the positive interaction with *recA*, *ssb* and *umuDC*.

A summary of other known interactions in the literature has been compiled by [6] and shown as the values between brackets in Table 1. We use this list as the 'ground truth' and compare the results of the network identification methods with it.

The plot in Fig. 8 shows a performance comparison of various identified models. Several observations that can be made from the comparison are as follows:

1. Incorporating the estimated error covariance in the error assessment improves the performance of the method. Comparing Models A and D, we can see that using the estimated error covariance gives us a smaller model with fewer false positives and negatives.
2. By increasing the value of  $\beta$ , we emphasise more on obtaining a smaller model (fewer connections) and less on the accuracy. Observing the results from Models B, C, D and E, we can see that it leads to fewer false positives and more false negatives. The model is the fewest total error (Model D) is shown in Fig. 7 panel (a).
3. The models identified using our method results in fewer errors compared to that in [6]. In particular, by tuning the





**Figure 7** Identification results of the Escherichia coli SOS network

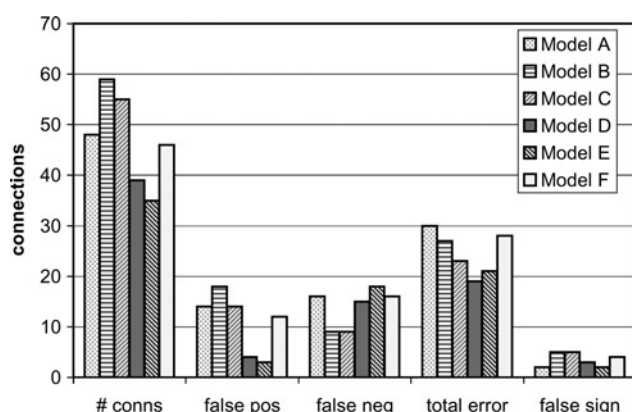
*a* The result of our method using estimated error covariance and  $\beta = 1.75$

*b* The network identified in [6]

Lines with arrows indicate activation, while circles indicate repression

Solid lines denote strong interaction, while dotted lines denote weak interaction

The distinction between strong and weak connections follows the same convention as in the segmentation polarity network



**Figure 8** Performance comparison of various identified models of the Escherichia coli SOS network

Model A is the result of our method without using error covariance estimation and  $\beta = 1.75$

Model B, C, D and E are the results of our method using error covariance estimation, with  $\beta = 1.1, 1.25, 1.75$  and  $2$ , respectively

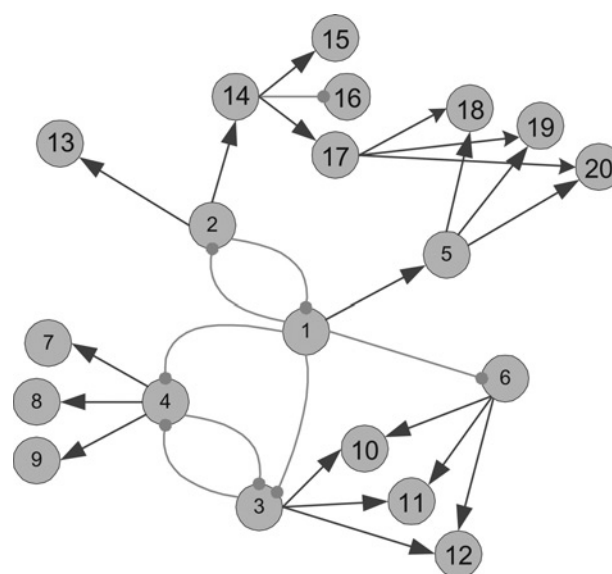
Model D is shown in Fig. 7 panel (a)

Model F is the network identified in [6], as also shown in Fig. 7 panel (b)

value of  $\beta$  carefully, we can obtain models with very low number of false positives (less than 5%).

### 4.3 Heterozygous knockdown of an in silico network

We generate an artificial network with 20 genes and study the performance of our method for various  $\beta$  parameters. As the measure of performance, we use the receiver operating



**Figure 9** The in silico network used in Section 4.3

characteristic (ROC) curve. The ROC curve plots the sensitivity of the prediction results against  $(1 - \text{specificity})$ . These quantities are given by the following formula [25]

$$\text{Sensitivity} = \frac{TP}{TP + FN}, \quad \text{Specificity} = \frac{TN}{TN + FP} \quad (31)$$

where T = True', F = False', P = Positives' and N = Negatives'.

As the true system, we use the network shown in Fig. 9. This network has a fan-like topology with a few regulator genes and a larger number of genes that are regulated directly or indirectly. The center of the network is formed by Genes 1 and 2, that are interconnected in a mutual inhibition. As such, these two genes form a toggle switch [26], in which they can only be active complementarily. Gene 1 acts as an inhibitor to another pair of genes (Genes 3 and 4) that also form a toggle switch. This group of four genes thus form a staged toggle switch than can hold one of the following three states: (on, off, off, off), (off, on, off, on) and (off, on, on, off). The remaining 16 genes in the network are regulated by these master genes directly or indirectly, as shown in Fig. 9.

We assume that in the wild-type specimen, each gene has two copies. Perturbation of the network is then performed by removing on the copies (heterozygous knockdown).

The dynamics of gene transcription and translation are modelled as a nonlinear system as follows

$$\frac{dx_i}{dt} = \left( \Gamma_0 + \prod_{j \in \text{Inh}_i} \frac{K^n}{K^n + y_j^n} \prod_{j \in \text{Act}_i} \frac{y_j^n}{K^n + y_j^n} \right) \Gamma_i - \lambda_i x_i \quad (32)$$

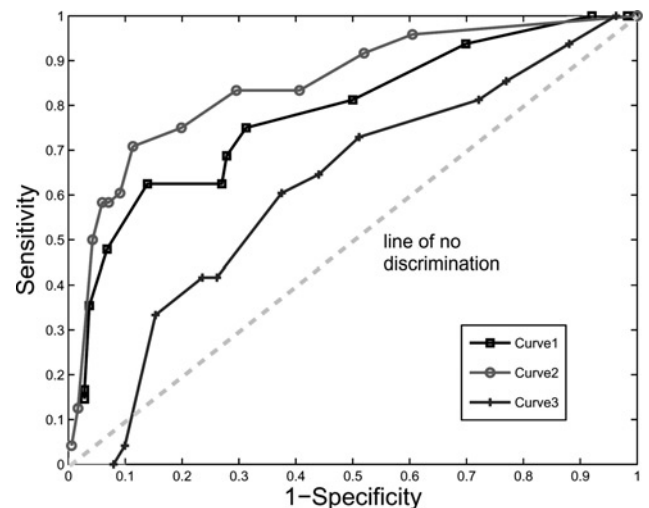
$$\frac{dy_i}{dt} = x_i - \kappa_i y_i, \quad i \in \{1, \dots, 20\} \quad (33)$$

where  $x_i$  is the transcript concentration of gene  $i$ ,  $y_i$  is the concentration of protein  $i$ ,  $\text{Inh}_i$  and  $\text{Act}_i$  are the set of genes that inhibit and activate gene  $i$ , respectively. The variable  $\Gamma_i$  represent the availability of gene  $i$ . In the wild type,  $\Gamma_{i,i \in \{1, \dots, 20\}} = 1$ . When gene  $i$  is knocked down,  $\Gamma_i = 1/2$ . The parameters of the model are  $\Gamma_0$  (basal transcription rate),  $K$  (inhibition and activation threshold),  $n$  (Hill coefficient),  $\lambda_i$  (transcript decay rate), and  $\kappa_i$  (protein decay rate).

We generate the data that we use in the identification by computing the steady state values of the transcript concentrations ( $x$ ) in (33), for the wild-type case and for each of the perturbations. The perturbation input  $U$  is taken to be a negative identity matrix of size 20. To simulate noisy measurement, we add some zero-mean Gaussian noise with uniform standard deviation to the computed transcript concentrations.

As we have seen in the previous sections, tuning the parameter  $\beta$  affects the performance of our algorithm. Since  $\beta$  represents a trade-off between sparsity and model completeness that cannot be known a priori, a fair assessment of the performance of the algorithm should be done by testing the algorithm for a wide range of choice of  $\beta$ . This is plotted in Fig. 10.

**4.3.1 Non-weighted recursive algorithm:** This algorithm is based on the approach taken in [21, 22]. This



**Figure 10** The ROC curves of various convex  $\ell_1$  relaxation techniques

Curve 1 is the result of non-weighted recursive algorithm, Curve 2 is the result of our algorithm and Curve 3 is the result of the non-recursive regularisation algorithm

algorithm does not use any weighting scheme in the  $\ell_1$  optimisation. Instead, in each iteration of the sparsity optimisation, entries that are small ( $< \delta$ ) are constrained to be zero in the subsequent iterations. The recursion is performed until it converges. The value of  $\delta$  can be tuned to adjust the sparsity of the result.

#### 4.3.2 Non-recursive regularisation algorithm:

This algorithm based on the approach taken in, for example [18–20]. As suggested by the name, this algorithm does not involve any iteration. To obtain a sparse model, a term with the  $\ell_1$  norm of the identified model is added to the cost function. Therefore instead of minimising  $\sum_{j=1, \dots, m} \eta_j^T R^j \eta_j$  in (24), we minimise  $\sum_{j=1, \dots, m} \eta_j^T R^j \eta_j + w \cdot \sum_{i,j} |A_{i,j}|$ . The weight  $w$  can be tuned to adjust the sparsity of the result. Specifically, larger  $w$  means higher priority on the sparsity of the model, and thus results in a sparser model.

The performance comparison between our method and these two methods are plotted in Fig. 10. We can see that for this example, our algorithm performs better than the other two relaxation techniques.

## 5 Conclusion

We propose a method for identifying genetic regulatory networks using expression profiles from genetic perturbation experiments. Some of the features of our method are, first, we aim at deriving a minimal model (characterised by the least number of connections) that explains the experimental data. Second, we can incorporate a priori information about the structure of the network. Third, we take into account the statistics of the measurement noise in formulating the cost function of our identification optimisation.

Our method is based on convex programming relaxation, that approaches the combinatorially hard problem of finding a minimal model with efficient computational scheme. In this paper, we test our method in a prototypical implementation that handles module size networks. However, as efficient customised implementation of convex optimisation algorithms are known to handle problems with thousands of variables, our method has a potential of solving the identification problem on a much larger scale.

## 6 Acknowledgment

The authors would like to thank Adam Halász, Marcin Imielinski, and Harvey Rubin for valuable discussion during the preparation of this paper. This work is partially funded by the ARO MURI SWARMS grant W911NF-05-1-0219, the NSF award ECS-0423905, the NSF award 0529426, and the NASA award NNX07AEIIA.

## 7 References

- [1] BONNEAU R., REISS D.J., SHANNON P., ET AL.: 'The inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets *de novo*', *Genome Biol.*, 2006, **7**, p. R36
- [2] HAYETE B., GARDNER T.S., COLLINS J.J.: 'Size matters: network inference tackles the genome scale', *Mole. Syst. Biol.*, 2007, **3**, article no. 77
- [3] GEIER F., TIMMER J., FLECK C.: 'Reconstructing gene-regulatory networks from time-series, knockout data and prior knowledge', *BMC Syst. Biol.*, 2007, **1**, (11), article no. 11
- [4] FAITH J.J., HAYETE B., THADEN J.T., ET AL.: 'Large-scale mapping and validation of *Escherichia coli* transcriptional regulation from a compendium of expression profiles', *PLoS Biol.*, 2007, **5**, (1), p. e8
- [5] BANSAL M., BELCASTRO V., AMBESI-IMPIOMBATO A., DI BERNARDO D.: 'How to infer gene networks from expression profiles', *Mole. Syst. Biol.*, 2007, **3**, article no. 78
- [6] GARDNER T.S., DI BERNARDO D., LORENZ D., COLLINS J.J.: 'Inferring genetic networks and identifying compound mode of action via expression profiling', *Science*, 2003, **301**, pp. 102–105
- [7] TEGNER J., YEUNG M.K.S., HASTY J., COLLINS J.J.: 'Reverse engineering gene networks: integrating genetic perturbations with dynamical modeling', *Proc. Natl. Acad. Sci.*, 2003, **100**, (10), pp. 5944–5949
- [8] SONTAG E., KIYATKIN A., KHOLODENKO B.N.: 'Inferring dynamic architecture of cellular networks using time series of gene expression, protein and metabolite data', *Bioinformatics*, 2004, **20**, (12), pp. 1877–1886
- [9] BANSAL M., DELLA GATTA G., DI BERNARDO D.: 'Inference of gene regulatory networks and compound mode of action from time course gene expression profiles', *Bioinformatics*, 2006, **22**, (7), pp. 815–822
- [10] HOFFMAN A.J., MCCORMICK S.T.: 'A fast algorithm that makes matrices optimally sparse', in PULLEYBANK W.R. (ED.): *Progress in combinatorial optimization*, (Academic Press, 1984), pp. 185–196
- [11] YEUNG M.K.S., TEGNER J., COLLINS J.J.: 'Reverse engineering gene networks using singular value decomposition and robust regression', *Proc. Natl. Acad. Sci.*, 2002, **99**, (9), pp. 6163–6168
- [12] THOMAS R., MEHROTRA S., PAPOUTSAKIS E.T., HATZIMANIKATIS V.: 'A model-based optimization framework for the inference on gene regulatory networks from DNA array data', *Bioinformatics*, 2004, **20**, (17), pp. 3221–3235
- [13] ANDREC M., KHOLODENKO B.N., LEVY R.M., SONTAG E.: 'Inference of signaling and gene regulatory networks by steady-state perturbation experiments: structure and accuracy', *J. Theor. Biol.*, 2005, **232**, (3), pp. 427–441
- [14] BOYD S., VANDENBERGHE L.: 'Convex optimization' (Cambridge University Press, 2004), Available online at [www.stanford.edu/boyd/cvxbook/](http://www.stanford.edu/boyd/cvxbook/)
- [15] BOYD S.: ' $\ell_1$ -norm methods for convex cardinality problems', 2007, Available online at [www.stanford.edu/class/ee364b/](http://www.stanford.edu/class/ee364b/), Lecture Notes for EE364b (Stanford University)
- [16] LOBO M.S., FAZEL M., BOYD S.: 'Portfolio optimization with linear and fixed transaction costs', *Ann. Oper. Res.*, 2007, **152**, (1), pp. 376–394
- [17] HASSIBI A., HOW J., BOYD S.: 'Low-authority controller design via convex optimization', *AIAA J. Guid. Control Dyn.*, 1999, **22**, (6), pp. 862–872
- [18] LI F., YANG Y.: 'Recovering genetic regulatory networks from micro-array data and location analysis data', *Genome Inf.*, 2004, **15**, (2), pp. 131–140
- [19] HAN S., YOON Y., CHO K.H.: 'Inferring biomolecular interaction networks based on convex optimization', *Comput. Biol. Chem.*, 2007, **31**, (5–6), pp. 347–354
- [20] GUO Y., SCHUURMANS D.: 'Learning gene regulatory networks via globally regularized risk minimization', 'Comparative genomics' (Springer Verlag, Berlin, 2007), pp. 83–95

- [21] PAPACHRISTODOULOU A., RECHT B.: 'Determining Interconnections in chemical reaction networks'. Proc. American Control Conf., New York, USA, 2007, pp. 4872–4877
- [22] COSENTINO C., CURATOLA W., MONTEFUSCO F., BANSAL M., DI BERNARDO D., AMATO F.: 'Linear matrix inequalities approach to reconstruction of biological networks', *IET. Syst. Biol.*, 2007, **1**, (3), pp. 164–173
- [23] BOYD S., GRANT M.C.: 'CVX – MATLAB software for disciplined convex programming', 2005, <http://www.stanford.edu/~boyd/cvx/>
- [24] LUSETTI S.L., VOLOSHIN O.N., INMAN R.B., CAMERINI-OTERO R.D., COX M.M.: 'The DinI Protein Stabilizes RecA Protein Filaments', *J. Biol. Chem.*, 2004, **279**, (29), pp. 30037–30046
- [25] DE MUTH J.E.: 'Basic statistics and pharmaceutical statistical applications' (Chapman & Hall/CRC, 2006, 2nd edn.)
- [26] GARDNER T.S., CANTOR C.R., COLLINS J.J.: 'Construction of a genetic toggle switch in Escherichia coli', *Nature*, 2000, **403**, pp. 339–342