

On reverse engineering of gene interaction networks using time course data with repeated measurements

SUPPLEMENTARY MATERIAL

Morrissey, E.R., Juárez, M.A.

Systems Biology Centre

University of Warwick, Coventry CV4 7AL, UK

K. J. Denby

Warwick HRI, Systems Biology Centre

University of Warwick, Coventry CV4 7AL, UK

Burroughs, N.J.

Systems Biology Centre,

University of Warwick, Coventry CV4 7AL, UK

These supplementary materials are presented in order of appearance within the main paper. For completeness, we briefly recall the model specification.

1S THE MODEL

Denote by y_g^t the expression level of gene $g = 1, \dots, G$, measured at time $t = 1, \dots, T$. We model the interaction network as a linear AR(1) process,

$$y_g^{t+1} = \mu_g + \sum_{j=1}^G \tilde{\beta}_{jg} y_j^t + \varepsilon_g^t, \quad (1)$$

where μ_g is the basal expression level of gene g ; $\tilde{\beta}_{jg} = \gamma_{jg} \beta_{jg}$ measures the influence of gene j on gene g , with $\beta_{jg} \in \mathbb{R}$ and $\gamma_{jg} = 1$ if j regulates g and $\gamma_{jg} = 0$ otherwise; finally, ε_g^t is an idiosyncratic error term, centred at zero and with precision parameter λ_g , typically assumed to be Gaussian. We augmented the model with the parenthood (link) indicator variables $\Gamma = \{\gamma_{jg}\}$ which will be the basis for estimating the network topology.

In order to account for the additional uncertainty when having repeated measurements we assume that the regulation process can be captured by (1), but instead of actually observing y_g^t , we have noisy measurements, x_{gr}^t , such that

$$x_{gr}^t = y_g^t + \eta_{gr}^t, \quad r = 1 \dots, R, \quad (2)$$

with η_{gr}^t a zero mean error measurement term, with precision parameter τ_g , independent for all t, g, r . This error term is frequently assumed Gaussian; however, given that the measurement process can potentially produce outliers, we will use a Student- t specification, $\text{St}(\eta_{gr}^t \mid 0, \tau_g, \nu)$, such that $\text{Var}[\eta_{gr}^t] = \nu \tau_g^{-1} / (\nu - 2)$ provided the degrees of freedom, $\nu > 2$.

2S THE PRIOR

The likelihood for the Student- t measurement AR(1) model is,

$$\ell(\Theta; X) = \prod_{g=1}^G \prod_{t=1}^T \prod_{r=1}^R N(y_g^{t+1} | \mu_g + \tilde{\beta}_g y^t, \lambda_g) N(x_{gr}^t | y_g^t, \omega_{gr}^t \tau_g) \text{Ga}(\omega_{gr}^t | \nu/2, \nu/2). \quad (3)$$

Where $Y = \{y_g^t\}$ are the unobserved expression levels, $X = \{x_{gr}^t\}$ denote their surrogate measurements and $\Theta = \{\mu, B, \Gamma, \lambda, \tau, \nu\}$ collects all the parameters involved, with $\mu = \{\mu_1, \dots, \mu_G\}$; $B = \{\beta'_1, \dots, \beta'_G\} \in \mathbb{R}^{G \times G}$ and $\beta_g = \{\beta_{1g}, \dots, \beta_{Gg}\}$; $\Gamma = \{\gamma_{ij}\}$; $\lambda = \{\lambda_1, \dots, \lambda_G\}$; and $\tau = \{\tau_1, \dots, \tau_G\}$.

We specify a product form (independent) prior,

$$\pi(\Theta) = \pi(\rho) \pi(\nu) \left[\prod_{g=1}^G \pi(\mu_g) \pi(\beta_g) \pi(\lambda_g) \pi(\tau_g) \pi(\gamma_g) \right], \quad (4)$$

where,

$$\pi(\mu_g) = N(\mu_g | 0, k), \quad (5)$$

$$\pi(\beta_g) = N_G(\beta_g | \mathbf{0}, k_\beta I), \quad g = 1, \dots, G, \quad (6)$$

$$\pi(\lambda_g) = \text{Ga}(\lambda_g | a_\lambda, b_\lambda), \quad (7)$$

$$\pi(\tau_g) = \text{Ga}(\tau_g | a_\tau, b_\tau), \quad (8)$$

$$\pi(\gamma_g | \rho) = \prod_{j=1}^g \text{Ber}(\gamma_{jg} | \rho), \quad g = 1, \dots, G, \quad (9)$$

$$\pi(\rho) = \text{Be}(\rho | a_\rho, b_\rho), \quad (10)$$

$$\pi(\nu) = \text{Ga}(\nu | a_\nu, b_\nu). \quad (11)$$

Given that the data is standardised before performing the estimation (zero mean and unitary standard deviation for each time series), we set $k_\mu = k_\beta = 1/4$; i.e. the prior variance of any component of μ and B is four. In our experience, this is typically not over-informative for microarray data.

As mentioned in the paper, when repeated measurements are available it is easier to estimate τ than λ . Thus, we set $\{a_\tau, b_\tau\} = \{2, 1/100\}$ which renders a rather flat prior with mode at 100 and variance of 20000.

For the autoregressive precision λ , we used $\{a_\lambda, b_\lambda\} = \{1/10, 1/10\}$. Thus setting the prior mean at one and the variance at 10. The mode now does not exist.

Derived from the conditions given in the paper: $P[\nu \leq 30] \approx 0.6$ and $\text{Mode}[\nu] = 15$, it is straightforward to verify that $\{a_\nu, b_\nu\} = \{3.5, 0.15\}$.

In the absence of any prior information, we treat ρ as the probability of any given link to be present and thus use the corresponding reference prior, $\text{Be}(\rho | 1/2, 1/2)$ (Bernardo and Smith, 1994, p. 315).

3S THE SAMPLER

AR(1) Precisions The full conditional of λ_g , $g = 1, \dots, G$ is given by

$$\pi(\lambda_g | \rightarrow) \propto \lambda_g^{T/2+a\lambda-1} \exp \left[-\lambda_g \left(b\lambda + \frac{1}{2} \left(y_g^{t+1} - \mu_g - y^t \tilde{\beta}_g \right)' \left(y_g^{t+1} - \mu_g - y^t \tilde{\beta}_g \right) \right) \right]$$

and thus can be sampled from a gamma distribution.

Constant term μ_g is conditionally Gaussian, with mean and precision

$$m_g = \frac{\bar{y}_g^{t+1} - \bar{y}^t \tilde{\beta}_g}{\lambda_g + k_\mu/T} \quad \text{and} \quad \tau'_\mu = k_\mu + T \lambda_g,$$

respectively, where $\bar{y}_g^{t+1} = T^{-1} \sum_t y_g^{t+1}$ and $\bar{y} = T^{-1} \sum_t y^t$.

Connectivity The overall connectivity, ρ , is sampled from a $\text{Be}(\rho | S + a_\rho, G^2 + b_\rho - S)$, with $S = \sum_{i,j=1}^G \gamma_{ji}$.

Measurement precision For each gene $g = 1, \dots, G$, the measurement precision, τ_g are updated from a gamma distribution $\text{Ga}(\tau_g | a'_\tau, b'_\tau)$ with

$$a'_\tau = RT/2 + a_\tau \quad \text{and} \quad b'_\tau = b_\tau + \frac{1}{2} \sum_{t=1}^T \sum_{r=1}^R \omega_{gr}^t (x_{gr}^t - y_{gr}^t)^2.$$

Degrees of freedom We use a Metropolis-within-Gibbs strategy to draw a new value, $v^{(m)}$, with a gamma proposal with its mean fixed at the previous draw, $v^{(m-1)}$. We control for the acceptance rate to lie around 1/3 by tuning the proposal's coefficient of variation, cv . Thus, we propose a new $v^{(m)}$ from

$$\text{Ga}(v^{(m)} | cv^{-2}, cv^{-2}/v^{(m-1)}).$$

Coefficients and link probabilities The update of each indicator variable γ_{jg} is performed jointly with all the corresponding coefficients

$$\beta : \beta^a \rightarrow \beta^b \quad \text{and} \quad \gamma : 0 \rightarrow 1$$

with acceptance probability

$$\alpha = \min \left\{ \frac{\pi(\tilde{\beta}^b)}{\pi(\tilde{\beta}^a)} \frac{q(\beta^a | \gamma^a) q(\gamma^a)}{q(\beta^b | \gamma^b) q(\gamma^b)}, 1 \right\},$$

where the subscripts have been removed for clarity. Given that we propose γ symmetrically, $q(\gamma^a)/q(\gamma^b) = 1$. The Hastings ratio is then

$$\frac{q(\beta^a | \gamma^a)}{q(\beta^b | \gamma^b)} = \frac{\rho}{1-\rho} k_\beta^{1/2} \frac{\exp \left[\frac{1}{2} \mu_\beta^b \Sigma_\beta^{-1b} \mu_\beta^b \right] \left| \Sigma_\beta^b \right|^{1/2}}{\exp \left[\frac{1}{2} \mu_\beta^a \Sigma_\beta^{-1a} \mu_\beta^a \right] \left| \Sigma_\beta^a \right|^{1/2}}.$$

with Σ . the covariance matrix obtained by considering only the relevant gene expression vectors. For the opposite move *i.e.* switching a link off, we use the reciprocal of the ratio above.

Non-observables These are drawn from a Gaussian distribution, $N(y_g^t | m_g^t, p_g^t)$, with location

$$m_g^t = \frac{\lambda_g m_{\text{AR}} + \tau_g m_{\text{meas}}}{p_g^t} \quad \text{and precision} \quad p_g^t = \lambda_g (1 + \tilde{\beta}_{gg}^2) + \tau_g \sum_{r=1}^R \omega_{gr}^t,$$

where

$$m_{\text{AR}} = \sum_{i \neq g} \tilde{\beta}_{ig} (y_i^{t-1} - \tilde{\beta}_{gg} y_i^t) + \tilde{\beta}_{gg} (y_g^{t-1} + y_g^{t+1}) \quad \text{and} \quad m_{\text{meas}} = \sum_{r=1}^R \omega_{gr}^t x_{gr}^t.$$

4S DATA SETS

Time traces of the data sets used in Section 4 of the paper. The linear *in silico* data, Figure-S 1a, comprises 16 genes measured at 41 time points. The ODE data has five genes and 50 measurements in time, Figure-S 1b.

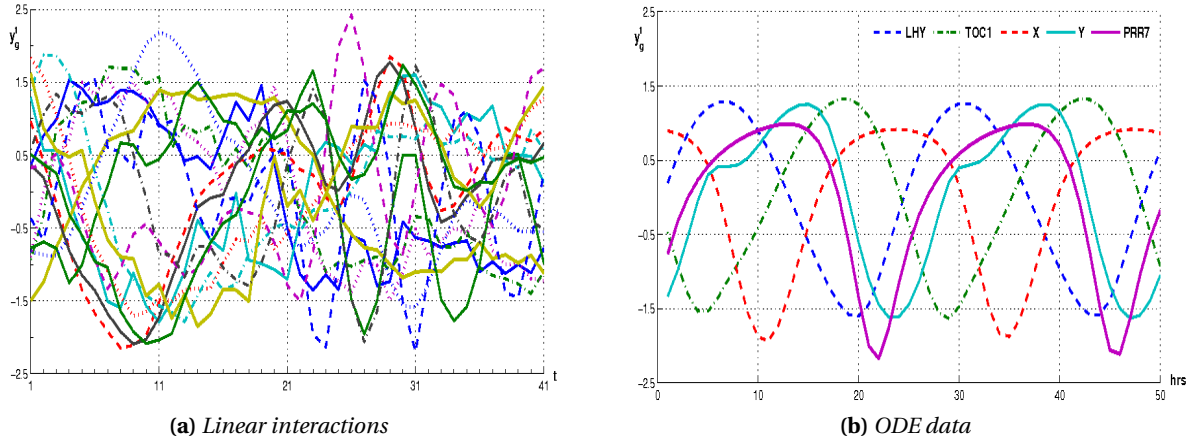


Figure-S 1. In silico data sets. Traces of the noiseless synthetic linear and ODE, non-linear data sets.

The *Arabidopsis* data set has 5 genes with 24 time points and four repetitions.

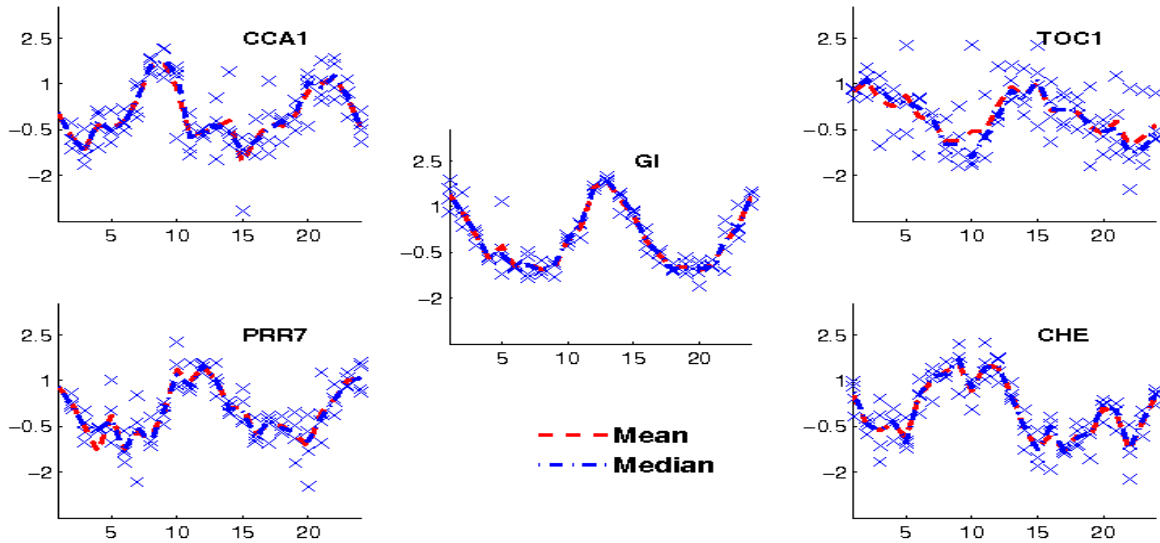


Figure-S 2. Circadian clock related genes in *Arabidopsis thaliana*. Gene expression repeated measurements. The mean (dashed) and median (dot-dashed) of each time point are plotted as time series.

4.1S MEAN CROSS ENTROPY

The score used in the paper, MxE, is simply the Kullback-Liebler divergency from any specific link to the true network configuration, $KL(\hat{y}_{ij} \mid p_{ij})$,

$$KL(\hat{y}_{ij} \mid p_{ij}) = p_{ij} \log \frac{p_{ij}}{\hat{y}_{ij}} + (1 - p_{ij}) \log \frac{1 - p_{ij}}{1 - \hat{y}_{ij}}$$

averaged over all possible links, $i, j = 1, \dots, G$, where $p_{ij} = 1$ if the link is present and zero otherwise; with the convention of $0 \log 0 = 0$.

4.2S *In silico* DATA

As expected, when a large number of replicates are available the three models yield similar networks, for a given threshold. However, the effect of overestimation in the regression coefficients posterior precisions is apparent when looking at the probabilities predicted by MM: again, they are less disperse than those predicted by either GM or SM.

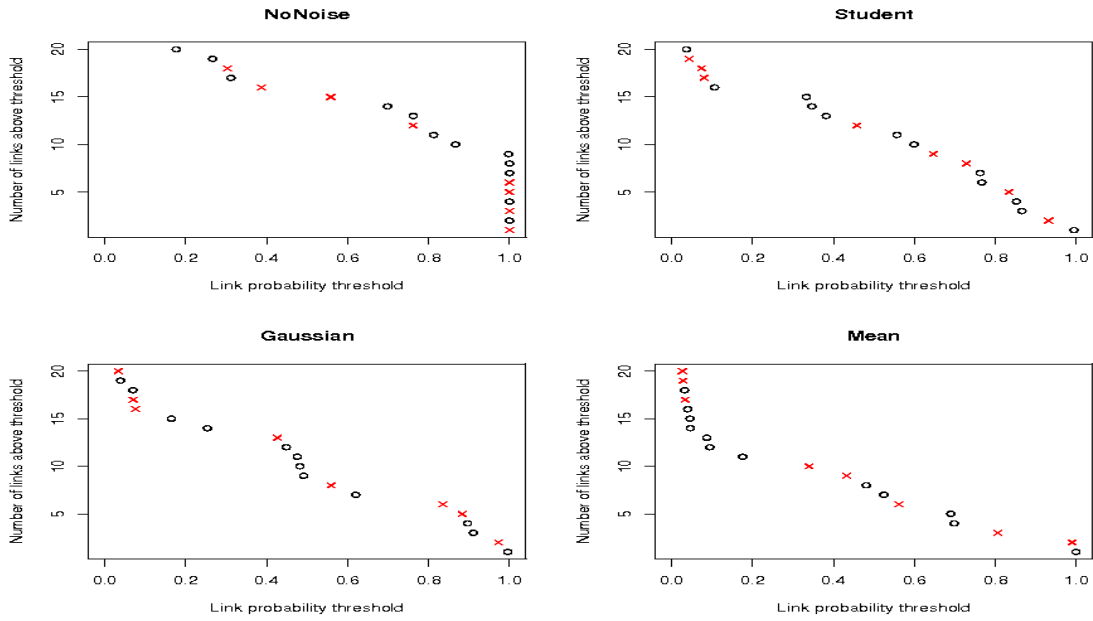


Figure-S 3. *ODE circadian clock in silico data set with $R = 20$ replicates. Number of links predicted present in the network versus posterior link probabilities for each model considered. A link present in the ODE model is highlighted with a cross.*

4.3S *In vivo* DATA

To have a more or less representative sample, we calculated all the possible subsamples with three replicates, \mathcal{P}_3 , and then measured the Euclidean distance between the standard deviations of the original data and \mathcal{P}_3 . These were classified into large, medium and small, based on their empirical distribution and ten series were selected from each region. For the 1-replicate case we used the Euclidean distance between the mean of the four replicate data set and the single data set. Interestingly we found no apparent effect of the Euclidean distance in the retrieved topologies and therefore we joined them when calculating the counts tables.

REFERENCES

Bernardo, J.M. and Smith, A.F.M. (1994) *Bayesian Theory*. Chichester: John Wiley & Sons.