Large data sets and complex models: A view from Systems Biology

Bärbel Finkenstädt, Department of Statistics, University of Warwick

OxWaSP workshop October 2015



Joint work with

Mathematical and Statistical Modelling:

Sylvia Calderazzo, Simone Tiberi, Kirsty Hey, Hiroshi Momiji, Dafyd Jenkins, George Minas, David Rand (University of Warwick)

Experimental Collaboration:

PRL expression in pituitary gland (Julien Davis, Mike White, University of Manchester)

Chronotherapy and Cancer Research Unit (Francis Lévi and group), Warwick Medical School and INSERM, Paris

Molecular Neurobiology of circadian timing (Michael Hastings and group), MRC Laboratory of Molecular Biology, Cambridge

PRESTA Consortium (University of Warwick)

Funding: BBSRC, EPSRC, MRC, EU (BIOSIM), Wellcome

Gene Expression



©1999 Addison Wesley Longman, Inc.

Standard Model of Gene Expression (Single Gene)



Standard Model of Gene Expression (Single Gene, ODE Version)

$$\frac{dm(t)}{dt} = \beta(t) - \delta_m m(t),$$

$$\frac{dp(t)}{dt} = \alpha m(t) - \delta_p p(t)$$

where

 $\boldsymbol{m}(t):$ concentration of mRNA

p(t): concentration of protein

 $\beta(t)$: transcription rate of mRNA

 δ_m, δ_p : degradation rate of mRNA, protein

Standard Model of gene expression (Stochastic version, single gene)

$$X(t) = \left(\begin{array}{c} X_m(t) \\ X_p(t) \end{array}\right)$$

4 reactions (transcription, degradation mRNA, translation, degradation protein)

$$v_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, v_2 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, v_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, v_4 = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$$

A reaction of type j changes X(t) to $X(t) + v_j$. Each reaction occurs at a rate $w_j(X(t))$.

Event	Effect	Transition Rate
Transcription	$(X_m, X_p) \to (X_m + 1, X_p)$	$w_1 = \beta(t)$
Degradation of mRNA	$(X_m, X_p) \to (X_m - 1, X_p)$	$w_2 = \delta_m X_m(t)$
Translation	$(X_m, X_p) \to (X_m, X_p + 1)$	$w_3 = \alpha X_m(t)$
Degradation of protein	$(X_m, X_p) \to (X_m, X_p - 1)$	$w_4 = \delta_p X_p(t)$

Table 1: Summary of reactions in the standard model of gene expression

Reaction networks constitute continuous time Markov jump processes and thus satisfy the Chapman-Kolmogorov equation for which one can obtain the forward form known as the master equation (ME) describing the evolution of the probability $P(X_m = n_1, X_p = n_2; t)$.

Although an exact numerical simulation algorithm is provided (Gillespie,1977) the ME is rarely tractable and hence an explicit formula for the exact likelihood is not available for parameter inference.



FIG. 2. Oscillations obtained by numerical simulation of the stochastic model for circadian rhythms with the Gillespie algorithm. The panels show oscillations of mRNA concentration, M (left column), limit cycles (second column), autocorrelation function (third column), and the period distribution (fourth column), for (A) $\Omega = 1000$, (B) $\Omega = 100$, and (C) $\Omega = 10$. The autocorrelation and period histograms have been calculated on a time series of 25 000 h, i.e., more than 1000 cycles. The white curve on the phase plane corresponds to the deterministic limit cycle. The deterministic oscillations have a period of 22 h. Parameter values are: $v_s = 1.6 \text{ nM h}^{-1}$, $K_I = 1 \text{ nM}$, n = 4, $v_m = 0.505 \text{ nM h}^{-1}$, $K_m = 0.5 \text{ nM}$, $k_s = 0.5 \text{ h}^{-1}$, $v_d = 1.4 \text{ nM h}^{-1}$, $K_d = 0.13 \text{ nM}$, $k_1 = 0.5 \text{ nM h}^{-1}$, $k_2 = 0.6 \text{ nM h}^{-1}$.

Gene Networks











ROR RBR loop REV-ERB, cytoplasm REV-ERB_N ROR. BMAL_ Ror Bmal Rev-Erb nucleus BMAL, ELOEK/BMAŁ Per Cry PER/ CRY PER*/ PER/ CRY CRY, PER*/ CRY_ PER/ CRY, $\operatorname{CRY}_{\operatorname{c}}$ PER_{c}^{*} PER_c PC loop

Figure 1. A model for the mammalian circadian clock.

Relógio A, Westermark PO, Wallach T, Schellenberg K, Kramer A, et al. (2011) Tuning the Mammalian Circadian Clock: Robust Synergy of Two Loops. PLoS Comput Biol 7(12): e1002309. doi:10.1371/journal.pcbi.1002309 http://127.0.0.1:8081/ploscompbiol/article?id=info:doi/10.1371/journal.pcbi.1002309





(8)

A

С

CLOCK/BMAL

 $\frac{dxI}{dt} = kf_{x1}x7 - kd_{x1}x1 - d_{x1}x1$

Rev-Erb



Ror



REV-ERB_C



ROR_C

$$\frac{dz7}{dt} = k_{p4}(y4 + y4_0) - ki_{z7}z7 - d_{z7}z7$$
(5)

REV-ERB_N

$$\frac{dx5}{dt} = ki_{z6}z6 - d_{x5}x5$$

ROR_N

 $\frac{dx6}{dt} = ki_{z7}z7 - d_{x6}x6$

BMAL_C

BMAL_N

Per

(1)

(2)

(3)

(4)

(6)

(7)

 $\frac{dz8}{dt} = k_{p5}(y5 + y5_0) - ki_{z8}z8 - d_{z8}z8$

$$\frac{dx7}{dt} = ki_{z8}z8 + kd_{x1}x1 - kf_{x1}x7 - d_{x7}x7$$

$$\frac{dy_1}{dt} = V_{1\max} \frac{1 + d\left(\frac{x_1}{k_{t_1}}\right)^b}{1 + \left(\frac{PC}{k_{t_1}}\right)^c \left(\frac{x_1}{k_{t_1}}\right)^b + \left(\frac{x_1}{k_{t_1}}\right)^b} - d_{y_1}y_1$$

CRY_C

$$\frac{dy2}{dt} = V_{2\max} \frac{1 + d\left(\frac{x1}{k_{t2}}\right)^e}{1 + \left(\frac{PC}{k_{t2}}\right)^f \left(\frac{x1}{k_{t2}}\right)^e + \left(\frac{x1}{k_{t2}}\right)^e} \frac{1}{1 + \left(\frac{x5}{k_{t21}}\right)^{f_1}} - d_{y_2}y_2$$

$$\frac{dz1}{dt} = k_{p2}(y2 + y2_0) + kd_{z4}z4 + kd_{z5}z5 - kf_{z5}z1z2 - kf_{z4}z1z3 - d_{z1}z^2$$





Gene Expression Data



Plant Microarray time series (PRESTA project)



Searching for transcription factors (TFs) of gene regulation from microarray data

Scenario:

Have (replicate) time series microarray gene expression data across various experiments and a set of candidate parents from Y1H experiments







Modeling:

Transcription $\beta(t)$ of a gene (*Child gene*) is regulated by transcriptional activators and/or inhibitors. These are protein products, called Transcription Factors (TFs), of other genes (Parent genes).

Wish to draw inference about regulation of mRNA transcription of a child gene by an unknown subset $\Gamma = \{\gamma_1, \gamma_2, \dots, \gamma_{\Gamma}\}$

of the candidate parents

$$\mathcal{G} = \{g_1, g_2, \dots, g_G\}$$

Switches occur at times where the expression, $P\gamma$ (t) of a parent $\gamma \in \Gamma$ crosses a threshold level due to an increase or decrease.

Define activation function

$$\alpha_{\Gamma}(t) = (\alpha_{\gamma_1}(t), \alpha_{\gamma_2}(t), \dots, \alpha_{\gamma_{\Gamma}}(t))$$

The set of switch times of the child's transcription is the set of time-points $\{s_1, s_2, \ldots, s_k\}$

at which at least one parent crosses its threshold.

The total time interval split into subintervals $[0, s_1], (s_1, s_2], \dots, (s_k, L]$ where the transcription rates of the subintervals are the same if the expression of all parent genes remains at the same state

For a given activation function $\alpha_{\Gamma}(t)$ we then have a set of parental states, $b_j, j = 1, 2, ..., \kappa$ observed in the union I_{b_j} of (at least one of) the above subintervals with the corresponding transcription rates $\tau_{b_j}, j = 1, 2, ..., \kappa$

$$\frac{dM}{dt} = \begin{cases} \tau_{b_1} - \delta_M M(t), & \text{for } t \in I_{b_1} \\ \tau_{b_2} - \delta_M M(t), & \text{for } t \in I_{b_2} \\ \vdots & \vdots \\ \tau_{b_\kappa} - \delta_M M(t), & \text{for } t \in I_{b_\kappa}. \end{cases}$$

Statistical Implementation:

Bayesian approach: RJMCMC to infer plausible models for parents

Likelihood based on normal error with inhomogeneous variance

Weighted Least Squares solution to piecewise linear ODE makes algorithm fast

Delayed RNA as proxy for unobserved TF's









Single cell imaging data



Reporter Gene constructs



Modelling single cell data

The class of state space models provides a unifying framework for modelling SRNs



h is the transition density of the approximating SRN

 $g\,$ is the density of the measurement process

Approximations

- ODE (neglects intrinsic noise)
- Chemical Langevin Equation (usually intractable)
- Linear Noise Approximation (linearization of the ME, tractable)
- Other approximations ?

Parameter Inference

The data likelihood is given by the marginal density

$$f(y|\theta) = \int_{x} f(y, x|\theta) dx$$

where the integrand can be factorized as

$$f(y, x|\theta) = h(x_0|\theta)g(y_0|x_0, \theta)\prod_{t=1}^T h(x_t|x_{t-1}, \theta)g(y_t|x_t, \theta)$$

Parameter Inference

Under the LNA (with Gaussian measurement error): Interval can be evaluated explicitly using Kalman methodology. Inference can be achieved by sampling from the posterior $f(\theta|y)$

Under BDA (for example) : use 2-step Gibbs sampler:

- 1. Sample the parameter vector from $f(\theta|x,y)$
- 2. Sample the latent states x from the filtering density $f(x|y,\theta)$







Challenges for statisticians in Systems Biology and Systems Medicine

- Complex and large networks of genes and their products
- Data from many sources:
- Replicates and if so what type?
- Many cells (Bayesian Hierarchical Modeling)
- Various labs (Prior Distributions)
- Destructive Sampling
-
- Different Types of Experiments imply different modeling approaches
- Information about intrinsic noise?
- Extrinsic noise?
- Destructive sampling?
- Reporter gene constructs? Which reporter gene?
- Type of camera used for imaging

- Not all state variables are observed Can use distributed delay functions as a proxy ?
- Parameter Identifiability
- Spatio-temporal modelling and inference
- Collaboration between experimentalists and mathematician/statisticians

[•]

[•]

Research Questions

- Mammalian Circadian Oscillator under diseases and treatments. Exploitation for Chronotherapy
- Development of signal processing tools for circadian time series to be used in forecasting systems for patient care at home
- Development of Spatio-temporal models to investigate how cells synchronize in the SCN

#B2 (-shBmal1)

Photon



Jour après inoculation tumorale

CT490 : KI/KI Per2::luc male mouse (B1)









ROR RBR loop REV-ERB, cytoplasm REV-ERB_N ROR. BMAL_ Ror Bmal Rev-Erb nucleus BMAL, ELOEK/BMAŁ Per Cry PER/ CRY PER*/ PER/ CRY CRY, PER*/ CRY_ PER/ CRY, $\operatorname{CRY}_{\operatorname{c}}$ PER_{c}^{*} PER_c PC loop

Figure 1. A model for the mammalian circadian clock.

Relógio A, Westermark PO, Wallach T, Schellenberg K, Kramer A, et al. (2011) Tuning the Mammalian Circadian Clock: Robust Synergy of Two Loops. PLoS Comput Biol 7(12): e1002309. doi:10.1371/journal.pcbi.1002309 http://127.0.0.1:8081/ploscompbiol/article?id=info:doi/10.1371/journal.pcbi.1002309

