

Parameter estimation for gene regulatory networks defined by differential equations

Nadia Lalam, Chalmers University of Technology

Joint work with Chris Klaassen, University of Amsterdam and EURANDOM

Outline

- Gene regulatory network
- Mathematical formalisms
- Statistical model
- Maximum likelihood based inference

Gene Regulatory Network

- Collection of genes in a cell which interact with each other and with other substances in the cell
- Regulation of gene expressions governs intracellular and extracellular mechanisms
- Aim: Determine the interactions within a gene regulatory network

Mathematical formalisms

- Many formalisms have been developed for qualitative or quantitative description of gene regulatory networks. They may be:
 - Discrete
 - Continuous
 - Stochastic
 - Deterministic

Mathematical formalisms

- Boolean Networks: Qualitative discrete models
 - State of a gene is either on or off
 - Connections between genes defining the activation/repression of a gene product on another
 - Change of state defined by an activation function
Example: $f_\ell(x) = \mathbf{1}(\sum_{k=1}^N W_{\ell k} x_k - h_\ell > 0)$,
 $x = (x_1, \dots, x_N) \in \{0, 1\}^N$
 - Synchronous or asynchronous

Mathematical formalisms

- Bayesian Networks: Probabilistic models
 - Vertices \leftrightarrow Random variables describing the gene expression levels
 - Conditional distributions of the vertices given their direct parents
 - Edges \leftrightarrow Dependencies
 - Static or dynamic

Mathematical formalisms

- Stochastic Equations:

- $P(\mathbf{X}, t)$ proba. that $\mathbf{X} = (X_1, \dots, X_N)$ molecules are inside the cell at time t , r reactions

$$P(\mathbf{X}, t + \delta t) = P(\mathbf{X}, t) \left(1 - \sum_{\ell=1}^r \alpha_{\ell} \delta t \right) + \sum_{\ell=1}^r \beta_{\ell} \delta t$$

$\alpha_{\ell} \delta t$ proba. that ℓ occurs during $[t, t + \delta t]$ given that \mathbf{X} molecules at time t , and $\beta_{\ell} \delta t$ proba. that the system is one ℓ reaction removed from the state \mathbf{X} at time t and then undergoes ℓ in $[t, t + \delta t]$

- $\delta t \rightarrow 0$ entails the Stochastic Master Equation

$$\frac{\partial P(\mathbf{X}, t)}{\partial t} = \sum_{\ell=1}^r (\beta_{\ell} - \alpha_{\ell} P(\mathbf{X}, t))$$

ODE's formalism

- State variables \leftrightarrow Gene product concentrations
- The gene product concentrations $g_{k\ell}(t; \theta)$ are assumed to satisfy a set of Ordinary Differential Equations

$$\frac{dg_{k\ell}(t; \theta)}{dt} = \psi_{k\ell}(\mathbf{g}(t; \theta), \theta)$$

k the cell index, $1 \leq k \leq K$

ℓ the gene product (protein or mRNA) index, $1 \leq \ell \leq N$

$\mathbf{g}(t; \theta) = [g_{k\ell}(t; \theta)]_{k,\ell}$

$\mathbf{g}(t_0; \theta)$ initial condition

θ the unknown biological parameter of interest

Example: Reaction-diffusion model

- Linear array of cells, N interacting genes
- Let $\mathbf{W} = [W_{\ell\ell'}]_{\ell, \ell'}$ be the regulatory weight matrix s.t.
 - $W_{\ell\ell'} > 0$ if gene ℓ' activates the synthesis of ℓ
 - $W_{\ell\ell'} < 0$ if gene ℓ' represses ℓ
 - $W_{\ell\ell'} = 0$ if ℓ' and ℓ do not interact
- Let D_ℓ be the diffusion parameter of gene product ℓ
- Let $\phi_\ell(x) = x^2 / (x^2 + h_\ell)$ (Hill function)

$$\frac{dg_{k\ell}(t; \theta)}{dt} = \phi_\ell \left(\sum_{\ell'=1}^N W_{\ell\ell'} g_{k\ell'}(t; \theta) \right) + D_\ell [g_{k-1\ell}(t; \theta) - 2g_{k\ell}(t; \theta) + g_{k+1\ell}(t; \theta)]$$

- Biological parameter: $\theta = (\mathbf{W}, \mathbf{D}, \mathbf{h})$

ODE's formalism

- k cell index, ℓ gene product index, initial conditions $\mathbf{g}(t_0; \theta)$ and

$$\frac{dg_{k\ell}(t; \theta)}{dt} = \psi_{k\ell}(\mathbf{g}(t; \theta), \theta)$$

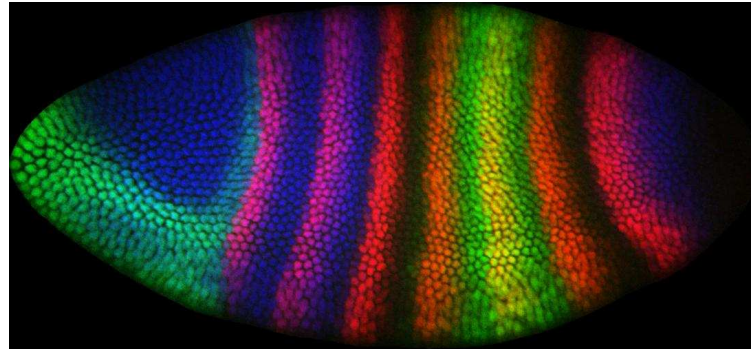
- Aim: Infer θ based on observations of the solutions of the ODE's at different times for different combinations of gene products

Proteomic data

- Scans of stained fixed organisms obtained by confocal laser scanning microscopy
- At maximum three stained gene products at a time whereas the gene regulatory network may contain more than three genes
 - ⇒ 2D picture
- Stained intensity proportional to the protein concentration
 - ⇒ Numerical data

Proteomic data

Image from confocal laser scanning microscopy...



... converted into numerical values

cell	caudal	hunchback	Kruppel
4113	33.57	91.41	9.28
4114	37.22	88.49	11.02
4115	44.17	83.57	12.91
4116	52.42	77.98	15.07
4117	59.76	71.47	17.71

...

Source: <http://flyex.ams.sunysb.edu/flyex/>

Poustelnikova et al., 2004

Statistical model

- Data obtained by confocal laser scanning microscopy divided into d samples: Each sample contains organisms stained for the same three gene products at the same intended time
- For an observation j from sample i , the random variable X_{ij} represents its scan converted into numerical values
- Aim: Infer θ , where $g_{k\ell}(t; \theta)$ are the solutions of the ODE's

Statistical model

- For a sample i formed by n_i observations, let L_i be the subset of three stained proteins among the N interacting genes, t_i the intended observation time
- X_{ij} is a table $\{X_{ijkl}\}_{k,\ell}$, where k indexes the cells and $\ell \in L_i$ indexes the gene products, such that

$$X_{ijkl} = g_{kl}(t_i + \delta_{ij}; \theta) + \varepsilon_{ijkl}$$
$$1 \leq i \leq d, 1 \leq j \leq n_i$$

- Account for stochasticity in the time of observation (random error δ_{ij}) and for stochasticity in the measurement of gene product (random error ε_{ijkl})

Statistical model

- Taylor approximation

$$g_{kl}(t_i + \delta_{ij}; \theta) \approx g_{kl}(t_i; \theta) + \frac{\partial g_{kl}(t; \theta)}{\partial t} \Big|_{t=t_i} \delta_{ij}$$

with $\partial g_{kl}(t_i; \theta) / \partial t = \psi_{kl}(\mathbf{g}(t_i; \theta), \theta) = \psi_{kl}(t_i; \theta)$ the RHS of the ODE at time $t = t_i$

- Then

$$X_{ijkl} = g_{kl}(t_i; \theta) + \psi_{kl}(t_i; \theta) \delta_{ij} + \varepsilon_{ijkl}$$

- Assume all ε_{ijkl} i.i.d. with density $f(\cdot)$, mean 0 and variance σ_ε^2 , all δ_{ij} i.i.d. with density $g(\cdot)$, mean 0 and variance σ_δ^2 , and ε_{ijkl} independent of δ_{ij}

Maximum Likelihood Estimation

- Assume $f(\cdot)$ and $g(\cdot)$ Gaussian densities
- Parameter $\gamma = (\theta, \sigma_\varepsilon^2, \sigma_\delta^2)$ estimated by MLE:
Maximization in γ of the likelihood

$$\prod_{i=1}^d \prod_{j=1}^{n_i} \int_{\mathbb{R}} \left[\prod_{k=1}^K \prod_{\ell \in L_i} \frac{1}{\sigma_\varepsilon} \phi \left(\frac{X_{ijkl} - g_{kl}(t_i; \theta) - \psi_{kl}(t_i; \theta) \sigma_\delta y}{\sigma_\varepsilon} \right) \right] \phi(y) dy$$

with $\phi(y) = 1/\sqrt{2\pi} \exp\{-0.5y^2\}$

Maximum Likelihood Estimation

- The log-likelihood normalized by $-2/n$, with $n = \sum_{i=1}^d n_i$ the number of observations, is

$$M_n(\gamma) = \frac{1}{n} \sum_{i,j} \left\{ K |L_i| \log \sigma_\varepsilon^2 + \log \left(1 + \frac{\sigma_\delta^2}{\sigma_\varepsilon^2} \sum_{k,l} \psi_{kl}^2(t_i; \theta) \right) \right. \\ \left. + \frac{1}{\sigma_\varepsilon^2} \sum_{k,l} (X_{ijkl} - g_{kl}(t_i; \theta))^2 \right. \\ \left. - \frac{\sigma_\delta^2 [\sum_{k,l} (X_{ijkl} - g_{kl}(t_i; \theta)) \psi_{kl}(t_i; \theta)]^2}{\sigma_\varepsilon^2 (\sigma_\varepsilon^2 + \sigma_\delta^2 \sum_{k,l} \psi_{kl}^2(t_i; \theta))} \right\}$$

with $\gamma = (\theta, \sigma_\varepsilon^2, \sigma_\delta^2)$

Maximum Likelihood Estimation

- Particular simple case: No time error ($\sigma_\delta = 0$). Then the MLE in the Gaussian case is the Least Squares Estimator
- Asymptotics: $n \rightarrow \infty$ s.t. $\lim_{n \rightarrow \infty} n_i/n = p_i > 0$
- The SLLN yields $\lim_{n \rightarrow \infty} M_n(\gamma) \stackrel{a.s.}{=} M(\gamma)$
- Let $\Gamma_0 = \{\gamma_* \in \Gamma : \gamma_* = \operatorname{argmin}_{\gamma \in \Gamma} M(\gamma)\}$ be the set of minimizers of the asymptotic criterion $M(\cdot)$
 $\Rightarrow \Gamma_0$ contains the true parameter value γ_0

Asymptotic properties

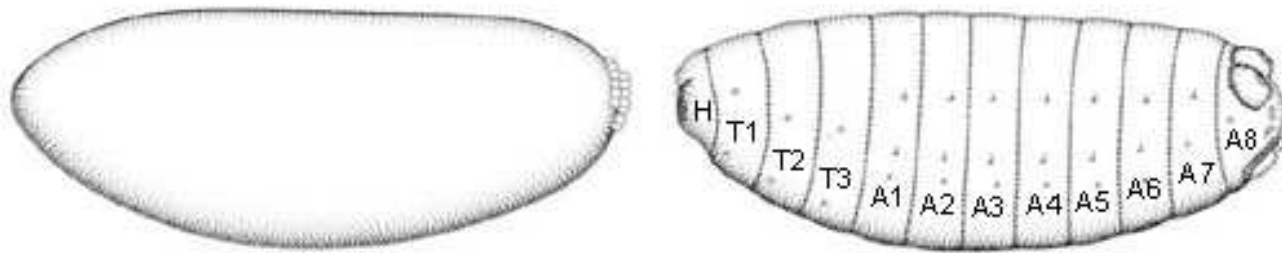
- $\Gamma_0 = \{\gamma_* \in \Gamma : \gamma_* = \operatorname{argmin}_{\gamma \in \Gamma} M(\gamma)\}$, and $n \rightarrow \infty$ s.t.
 $\lim_{n \rightarrow \infty} n_i/n = p_i > 0$
- Under smoothness conditions of the RHS of the ODE's, the MLE $\hat{\gamma}_n$ is consistent in the sense that, for all $\varepsilon > 0$, $\lim_{n \rightarrow \infty} P(\inf_{\gamma_* \in \Gamma_0} d(\hat{\gamma}_n, \gamma_*) \geq \varepsilon) = 0$
- If $\Gamma_0 = \{\gamma_0\}$, then $\hat{\gamma}_n$ is weakly consistent, i.e. for all $\varepsilon > 0$, $\lim_{n \rightarrow \infty} P(d(\hat{\gamma}_n, \gamma_0) \geq \varepsilon) = 0$

Asymptotic properties

- $\Gamma_0 = \{\gamma_* \in \Gamma : \gamma_* = \operatorname{argmin}_{\gamma \in \Gamma} M(\gamma)\}$, and $n \rightarrow \infty$ s.t.
 $\lim_{n \rightarrow \infty} n_i/n = p_i > 0$
- If $\hat{\gamma}_n$ is weakly consistent, and under regularity assumptions of the RHS of the ODE's and of the expectation of $M_n(\gamma)$, then $\hat{\gamma}_n$ is root- n consistent, that is $\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} P(\sqrt{n}d(\hat{\gamma}_n, \gamma_0) > M) = 0$

Drosophila case

- Work in progress: Application to the MLE approach to experimental data
- Gene regulatory network responsible of a particular early stage of Drosophila embryo development: Segmentation



Drosophila case

- $N = 6$ genes in the network responsible of the early embryo segmentation
- $n = 954$ observations between cleavage cycle 13 (after egg fertilization) and cleavage cycle 14
- Observations: Protein concentrations
- 9 time classes
- $K = 58$ cells along the embryo antero-posterior axis

Mjolsness et al., 1991, Jaeger et al., 2004

Drosophila case

- Initial conditions at the onset of cleavage cycle 13 and

$$\begin{aligned}\psi_{k\ell}(\mathbf{g}(t; \theta), \theta) &= R_\ell \Phi\left(\sum_{\ell'=1}^{N_g} W_{\ell\ell'} g_{k\ell'}(t) + m_\ell g_{k\ell}(t) + h_\ell\right) \\ &\quad + D_\ell [g_{k-1\ell}(t) - 2g_{k\ell}(t) + g_{k+1\ell}(t)] \\ &\quad - \lambda_\ell g_{k\ell}(t)\end{aligned}$$

with $\Phi(x) = 0.5[(x/\sqrt{x^2 + 1}) + 1]$

- Gene regulation and synthesis/Diffusion/Decay
- $\theta = ((R_\ell, m_\ell, h_\ell, D_\ell, \lambda_\ell)_{1 \leq \ell \leq N}, (W_{\ell\ell'})_{1 \leq \ell, \ell' \leq N})$

Mjolsness et al., 1991, Jaeger et al., 2004

Summary and Perspectives

● Summary

- Gene regulatory networks modelled with ODE's
- Statistical model for data corrupted by two sources of noise in state variable and time observation
- Inference of a parameter of ODE's by MLE
- Asymptotic properties and ongoing work to apply the MLE procedure to data

● Perspectives

- Remove the assumption of Gaussian errors
- Extend the approach to Partial Differential Equations
- Account for differences in shape between organisms

Bibliography

Bickel, P. J., Klaassen, C. A. J., Ritov, Y., Wellner, J. (1998) Springer

Bolouri, H., Davidson, E. H. (2002) *BioEssays*, 24, 1118-1129

Hasty, J., McMillen, D., Isaacs, F., Collins, J. J. (2001) *Nature Reviews Genetics*, 2, 268-279

Jaeger, J., Surkova, S., Blagov, M., Janssens, H., Kosman, D., Kozlov, K., Manu, Myasnikova, E., Vanario-Alonso, C. E., Samsonova, M., Sharp, D. H., Reinitz, J. (2004) *Nature*, 430, 368-371

de Jong, H. (2002) *Journal of Computational Biology*, 9, 67-103

Kosman, D., Small, S., Reinitz, J. (1998) *Development Genes and Evolution*, 208, 290-294

Lalam N., Klaassen C. (2006) EURANDOM Preprint 2006-018

Mjolsness, E., Sharp, D. H., Reinitz, J. (1991) *Journal of Theoretical Biology*, 152, 429-453

Poustelnikova, E., Pisarev, A., Blagov, M., Samsonova, M., Reinitz, J. (2004) *Bioinformatics*, 20, 2212-2221

Acknowledgements

- ▷ Joke Blom, CWI Amsterdam
- ▷ Richard Gill, University of Leiden and EURANDOM
- ▷ Jaap Kaandorp, University of Amsterdam
- ▷ Members of the Reinitz Fly Lab, Stony Brook University, New York