

CS-E5880 Modeling biological networks

Parameter estimation for biological networks

Harri Lähdesmäki

Acknowledgement: Jukka Intosalmi, Henrik Mannerström

Department of Computer Science
Aalto University

January 28, 2022

Outline

- ▶ ODE model calibration as an optimization problem
- ▶ Local gradient based optimization
- ▶ Sensitivity equations
- ▶ Practical considerations
- ▶ ODE model calibration as a statistical estimation problem
- ▶ Identifiability
- ▶ A collection of articles covering the topic listed at the end for additional reading

Parameter estimation

Parameter estimation is an important part of biological network modeling. Here we focus on ODE models but other models have parameters too.

Why parameter estimation?

- ▶ Estimation of the reaction rates and initial values
- ▶ Testing if a hypothetical model can produce observed dynamics
- ▶ Model construction in a data-driven manner
- ▶ Experimental data is typically noisy, so many parameter values can fit the data reasonably well

Parameter estimation as an optimization problem (1)

- ▶ Let us consider a one dimensional initial value problem

$$\frac{dx}{dt} = f(x, \theta), \quad x(0) = x_0,$$

where $\theta = (\theta_1, \dots, \theta_d)$ is a parameter vector.

- ▶ The system has the solution $x(t, \theta)$.
- ▶ The solution can be obtained numerically (recall Euler and Runge-Kutta methods)
- ▶ The observed values of $x(t, \theta)$ are denoted by

$$y(t_i), \quad i = 1, \dots, n$$

where t_i are the measurement times.

- ▶ Observed values can be corrupted by noise

Parameter estimation as an optimization problem (2)

- ▶ In optimization setting, our goal is to minimize a distance between the solution of the ODE model and data
- ▶ For example, we may wish to minimize the sum-of-squares objective function

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n (y(t_i) - x(t_i, \boldsymbol{\theta}))^2$$

- ▶ For ODE models, closed-form solutions are only rarely available and **numerical optimization techniques** need to be used.
- ▶ Typically a non-convex optimization problem

Optimization techniques

Global and local optimization techniques

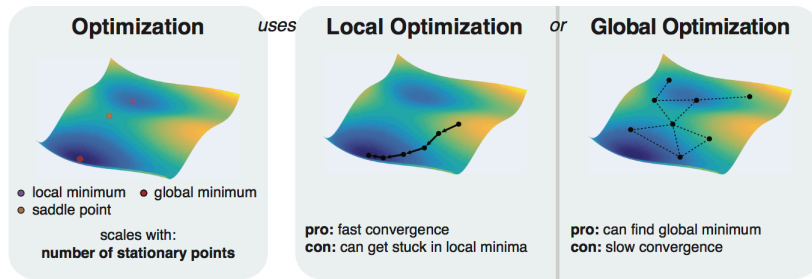
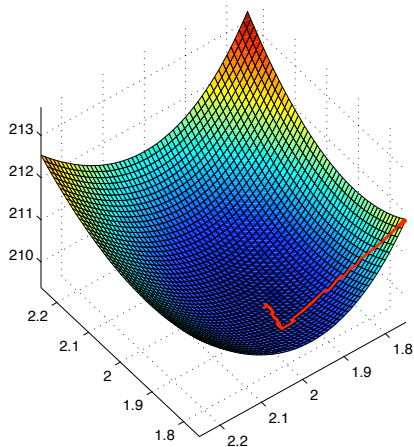


Figure: Illustration of optimization problem (from Fröchlich et al., 2017)

- ▶ In this lecture, we will concentrate on deterministic local optimization
 - ▶ Extensions to global optimization through multistart approach

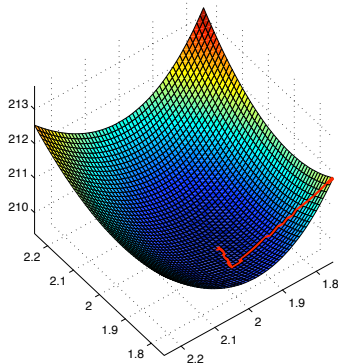
Gradient based optimization

- Search for the nearest minimum of the objective function following the direction of the gradient



Gradient descent method

1. Set $k = 0$ and define initial parameter values θ_k as well as step size (learning rate) r
2. Evaluate the objective function $L(\theta_k)$
3. Compute the gradient of the objective function $\nabla L(\theta) = \left(\frac{\partial L(\theta)}{\partial \theta_1}, \dots, \frac{\partial L(\theta)}{\partial \theta_d} \right)$ at θ_k , $\nabla L(\theta_k)$
- [4.] Optional line search: solve for the \hat{r} minimizing $L(\theta_k - r \cdot \nabla L(\theta_k))$, set $r := \hat{r}$
5. Set $\theta_{k+1} = \theta_k - r \cdot \nabla L(\theta_k)$
6. Compare the current objective function value with the previous value;
 - ▶ if $L(\theta_{k+1}) > L(\theta_k)$ or
 - ▶ if $\|\theta_{k+1} - \theta_k\| < \epsilon$then the solution is θ_k
7. Otherwise, set $k := k + 1$ and go to step 2



Differentiation of the objective function

- To obtain the gradient, we need to differentiate the objective function

$$\begin{aligned}\frac{\partial L(\boldsymbol{\theta})}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \sum_{i=1}^n (y(t_i) - x(t_i, \boldsymbol{\theta}))^2 \\ &= \sum_{i=1}^n -2(y(t_i) - x(t_i, \boldsymbol{\theta})) \frac{\partial x(t_i, \boldsymbol{\theta})}{\partial \theta_j}\end{aligned}$$

Differentiation of the objective function

- To obtain the gradient, we need to differentiate the objective function

$$\begin{aligned}\frac{\partial L(\boldsymbol{\theta})}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \sum_{i=1}^n (y(t_i) - x(t_i, \boldsymbol{\theta}))^2 \\ &= \sum_{i=1}^n -2(y(t_i) - x(t_i, \boldsymbol{\theta})) \frac{\partial x(t_i, \boldsymbol{\theta})}{\partial \theta_j}\end{aligned}$$

- The gradient depends on $\frac{\partial x(t_i, \boldsymbol{\theta})}{\partial \theta_j}$, $i = 1, \dots, n; j = 1, \dots, d$
- Note that $x(t_i, \boldsymbol{\theta})$ is a function (or infinite dimensional)
- Three approaches:
 1. Finite differences approximation
 2. Sensitivity equations (forward accumulation)
 3. Adjoint method (reverse accumulation)

Finite differences approximation for $\frac{\partial x(t_i, \boldsymbol{\theta})}{\partial \theta_j}$

- ▶ The approximation is computed by perturbing the parameters

$$\begin{aligned}\frac{\partial x(t_i, \boldsymbol{\theta})}{\partial \theta_j} &= \lim_{\varepsilon \rightarrow 0} \frac{x(t_i, \boldsymbol{\theta}) - x(t_i, \boldsymbol{\theta} + \varepsilon \mathbf{e}_j)}{\varepsilon} \\ &\approx \frac{x(t_i, \boldsymbol{\theta}) - x(t_i, \boldsymbol{\theta} + h \mathbf{e}_j)}{h}\end{aligned}$$

where \mathbf{e}_j is the j th unit vector and h is a sufficiently small constant.

- ▶ Easy to implement numerically, but...
 - ▶ Requires several numerical solutions for the ODE system
 - ▶ No generic way of choosing the constant h
 - ▶ In general finite differences approach may result in a poor approximation
 - ▶ Better results can be obtained using the sensitivity equations

Sensitivity equations

- ▶ We can take a so-called total derivative of the ODE system w.r.t. θ_j

$$\frac{d}{d\theta_j} \frac{dx}{dt} = \frac{d}{d\theta_j} f(x, \theta)$$

Sensitivity equations

- ▶ We can take a so-called total derivative of the ODE system w.r.t. θ_j

$$\frac{d}{d\theta_j} \frac{dx}{dt} = \frac{d}{d\theta_j} f(x, \theta)$$

- ▶ Recall that the derivative of a composite function $f(g(x))$ is $f'(g(x))g'(x)$
- ▶ Here x depends on θ (and time), $x(\theta)$, so

$$f(x(\theta), \theta) = f(\underbrace{(x(\theta), \theta)}_{\text{function } g(\theta)})$$

Sensitivity equations

- The total derivative is thus

$$\frac{d}{dt} \frac{dx}{d\theta_j} = \frac{d}{d\theta_j} f(x(\theta), \theta) = f'(\underbrace{(x(\theta), \theta)}_{\text{function } g(\theta)}) \underbrace{\left(\frac{dx(\theta)}{d\theta_j} \frac{d\theta}{d\theta_j} \right)^T}_{\text{derivative } g'(\theta)}$$

Sensitivity equations

- The total derivative is thus

$$\begin{aligned}\frac{d}{dt} \frac{dx}{d\theta_j} &= \frac{d}{d\theta_j} f(x(\theta), \theta) = f'(\underbrace{(x(\theta), \theta)}_{\text{function } g(\theta)}) \underbrace{\left(\frac{dx(\theta)}{d\theta_j} \quad \frac{d\theta}{d\theta_j} \right)^T}_{\text{derivative } g'(\theta)} \\&= \left(\frac{\partial f(x, \theta)}{\partial x} \quad \frac{\partial f(x, \theta)}{\partial \theta} \right) \left(\frac{dx}{d\theta_j} \quad \frac{d\theta}{d\theta_j} \right)^T \\&= \frac{\partial f(x, \theta)}{\partial x} \frac{dx}{d\theta_j} + \frac{\partial f(x, \theta)}{\partial \theta} \frac{d\theta}{d\theta_j} \\&= \frac{\partial f(x, \theta)}{\partial x} \frac{dx}{d\theta_j} + \frac{\partial f(x, \theta)}{\partial \theta_j}\end{aligned}$$

(note that we change the notation from $x(\theta)$ back to x after the first line)

Sensitivity equations

- ▶ The total derivative

$$\frac{d}{dt} \frac{dx}{d\theta_j} = \frac{\partial f(x, \theta)}{\partial x} \frac{dx}{d\theta_j} + \frac{\partial f(x, \theta)}{\partial \theta_j}$$

- ▶ Let us define the sensitivity w.r.t. the j th parameter as $s_j(t) = \frac{\partial x(t, \theta)}{\partial \theta_j}$
- ▶ The sensitivity equation can then be written as

$$\frac{ds_j}{dt} = \frac{\partial f(x, \theta)}{\partial x} s_j + \frac{\partial f(x, \theta)}{\partial \theta_j}$$

- ▶ This is another differential equation

Solving the sensitivities

- ▶ The sensitivities and the original ODE model form a coupled system and they are solved simultaneously using numerical solvers.
- ▶ The accuracy of the computed sensitivities depends only on the numerical error of the ODE solver!

Comparing the finite differences and sensitivity equation approaches

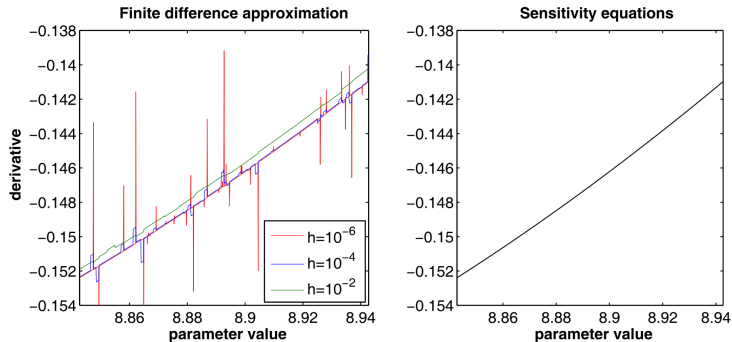


Figure: From Raue et al., 2013

Example of sensitivity equation construction (1)

- ▶ Let us consider the initial value problem

$$\frac{dx}{dt} = \underbrace{\alpha + \theta_1 \frac{x}{x+K} - \theta_2 x}_{f(x,\theta)}, \quad x(0) = x_0$$

where θ_1 and θ_2 are unknown rate parameters.

- ▶ We have

$$\begin{aligned}\frac{\partial f(x, \theta)}{\partial x} &= \theta_1 \frac{K}{(x+K)^2} - \theta_2 \\ \frac{\partial f(x, \theta)}{\partial \theta_1} &= \frac{x}{x+K} \\ \frac{\partial f(x, \theta)}{\partial \theta_2} &= -x.\end{aligned}$$

Example of sensitivity equation construction (2)

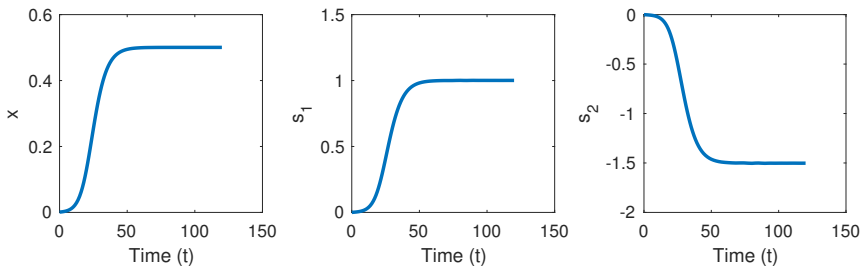
- ▶ The sensitivity equations together with the actual ODE model form the following ODE system

$$\begin{aligned}\frac{dx}{dt} &= \alpha + \theta_1 \frac{x}{x+K} - \theta_2 x \\ \frac{ds_1}{dt} &= \left(\theta_1 \frac{K}{(x+K)^2} - \theta_2 \right) s_1 + \frac{x}{x+K} \\ \frac{ds_2}{dt} &= \left(\theta_1 \frac{K}{(x+K)^2} - \theta_2 \right) s_2 - x\end{aligned}$$

with $x(0) = x_0 = 0$ and $s_1(0) = s_2(0) = 0$.

Example of sensitivity equation construction (3)

- And the system can be solved.



- The final gradient computation via the derivative of the objective function

$$\frac{\partial L(\theta)}{\partial \theta_j} = \sum_{i=1}^n -2(y(t_i) - x(t_i, \theta)) \frac{\partial x(t_i, \theta)}{\partial \theta_j}$$

Sensitivity equation in the general form

- ▶ Let us consider an m -dimensional initial value problem ($\mathbf{x} \in \mathbb{R}^m$, $\boldsymbol{\theta} \in \mathbb{R}^d$, $\mathbf{f}(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^d \rightarrow \mathbb{R}^m$)

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, \boldsymbol{\theta}), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)^T$ is a parameter vector.

- ▶ The sensitivity equations can be expressed in the general form

$$\underbrace{\frac{d}{dt} \frac{d\mathbf{x}}{d\boldsymbol{\theta}}}_{\in \mathbb{R}^{m \times d}} = \underbrace{\frac{\partial \mathbf{f}(\mathbf{x}, \boldsymbol{\theta})}{\partial \mathbf{x}}}_{\in \mathbb{R}^{m \times m}} \underbrace{\frac{d\mathbf{x}}{d\boldsymbol{\theta}}}_{\in \mathbb{R}^{m \times d}} + \underbrace{\frac{\partial \mathbf{f}(\mathbf{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}}_{\in \mathbb{R}^{m \times d}}$$

and solved together with the original ODE shown above

Adjoint method

- ▶ The original problem of solving $x(t, \theta)$ forward in time is transformed into an *adjoint problem* of solving $\lambda(t, \theta)$ backward in time.
- ▶ Integrate the original problem $x(t, \theta)$ forward from t_1 to t_n to obtain the trajectory $x(t_1, \theta), \dots, x(t_n, \theta)$.
- ▶ Integrate the adjoint problem $\lambda(t, \theta)$ backward from t_n to t_1 to obtain the trajectory $\lambda(t_1, \theta), \dots, \lambda(t_n, \theta)$.
- ▶ Extract the derivatives from the adjoint solution.

Practical considerations

- ▶ In practice, it can be beneficial to carry out the parameter estimation using log-transformed parameters
- ▶ Automated construction of the sensitivity equation and adjoint equations.
- ▶ Selection of the numerical solver (e.g. Sundials CVODES)
- ▶ Least squares optimization has been used in many applications

Global optimization through multistart approach

- ▶ A non-convex optimization problem
- ▶ Multiple local minima
- ▶ Parameter space explored using multiple optimization runs
- ▶ Parallelizable for high-dimensional problems

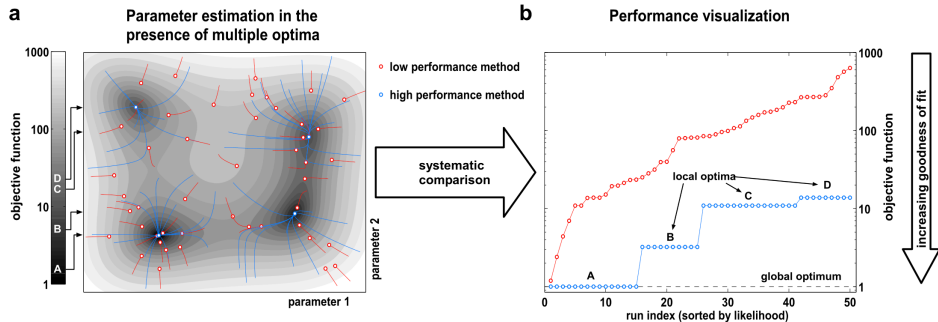


Figure: Illustration of multistart approach (from Raue et al, 2013)

Choosing the initial points for multistart (1)

- ▶ Random (Monte Carlo)
- ▶ Latin hypercube sampling (LHS)
 - ▶ Divide the range of each variable θ_j into equal number (equally probable) of bins
 - ▶ Sampling such that for each variable θ_j a particular bin is chosen only once
 - ▶ Thus, sampling does not depend on dimension (d) of θ

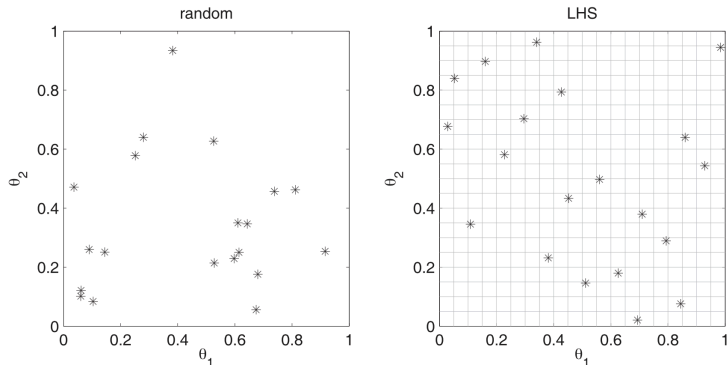


Figure: From Raue et al, 2013

Choosing the initial points for multistart (2)

- Latin hyper-cube sampling gives initial parameter values which are on average more far away from each other

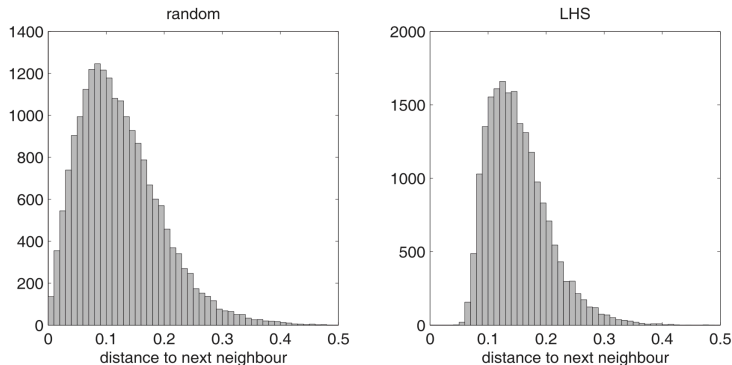


Figure: From Raue et al, 2013

Statistical formulation of the parameter estimation problem

- ▶ Formulating what is known about the measurement system in statistical terms
- ▶ If we assume normally distributed measurement errors, we can write

$$y(t_i) = x(t_i, \boldsymbol{\theta}) + \epsilon(t_i),$$

where $\epsilon(t_i) \sim \mathcal{N}(0, \sigma_i^2)$.

- ▶ This is equivalent to

$$y(t_i) \sim \mathcal{N}(x(t_i, \boldsymbol{\theta}), \sigma_i^2).$$

- ▶ By assuming independent measurements, we can express the likelihood of the data given the parameters in the form

$$l(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{2\sigma_i^2}\right).$$

Maximum likelihood estimation

- ▶ The maximum likelihood principle is to choose the value $\theta = \hat{\theta}$ which maximizes $l(\theta)$.
- ▶ In other words, the maximum likelihood estimate is

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

- ▶ It is equivalent to solve

$$\hat{\theta} = \arg \min_{\theta} L(\theta),$$

where $L(\theta) = -2 \log(l(\theta))$ is the negative log-likelihood.

Negative log-likelihood (1)

$$\begin{aligned} L(\boldsymbol{\theta}) &= -2 \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left(-\frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{2\sigma_i^2} \right) \right) \\ &= -2 \left(\sum_{i=1}^n -\frac{1}{2} \log(2\pi\sigma_i^2) - \frac{1}{2} \left(\frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{\sigma_i^2} \right) \right) \\ &= \sum_{i=1}^n \left(\log(2\pi) + 2 \log(\sigma_i) + \frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{\sigma_i^2} \right) \\ &= \sum_{i=1}^n \left(2 \log(\sigma_i) + \frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{\sigma_i^2} \right) + C_1 \end{aligned}$$

Negative log-likelihood (2)

- If $\sigma_i = \sigma, i = 1, \dots, n$, we have

$$\begin{aligned} L(\boldsymbol{\theta}) &= \sum_{i=1}^n \left(2 \log(\sigma) + \frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{\sigma^2} \right) + C_1 \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (y(t_i) - x(t_i, \boldsymbol{\theta}))^2 + C_2 \end{aligned}$$

Negative log-likelihood (2)

- ▶ If $\sigma_i = \sigma, i = 1, \dots, n$, we have

$$\begin{aligned} L(\boldsymbol{\theta}) &= \sum_{i=1}^n \left(2 \log(\sigma) + \frac{(y(t_i) - x(t_i, \boldsymbol{\theta}))^2}{\sigma^2} \right) + C_1 \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (y(t_i) - x(t_i, \boldsymbol{\theta}))^2 + C_2 \end{aligned}$$

- ▶ In general, measurement uncertainty can be characterized by any (possibly non-Gaussian) distribution

$$y(t_i) \sim g(x(t_i, \boldsymbol{\theta}) | \theta_{\text{meas}})$$

Parameter identifiability

- ▶ In many cases, the parameters are not well determined
- ▶ Infinitely many parameter settings may have equal likelihood
- ▶ The effect of changing one parameter can be compensated by tuning other parameters

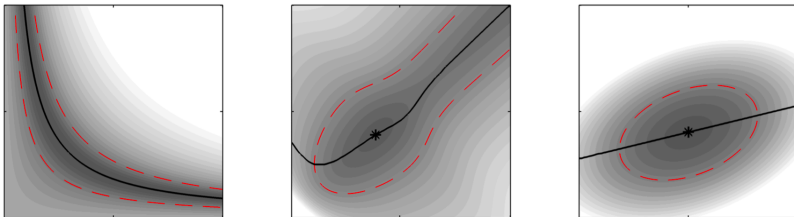


Figure: From Raue et al, 2010

Profile likelihood (1)

- ▶ Parameter identifiability can be studied e.g. using the profile likelihood method.
- ▶ The profile likelihood for the j th variable is defined by

$$PL_j(p) = \max_{\theta \in \{\theta | \theta_j = p\}} L(\theta),$$

where L is the log-likelihood.

- ▶ Evaluate the profile likelihood around the optimal parameter values θ^*
- ▶ Structural and practical non-identifiability
- ▶ Sometimes non-identifiability can be solved by reparameterization of the ODE system

Profile likelihood (2)

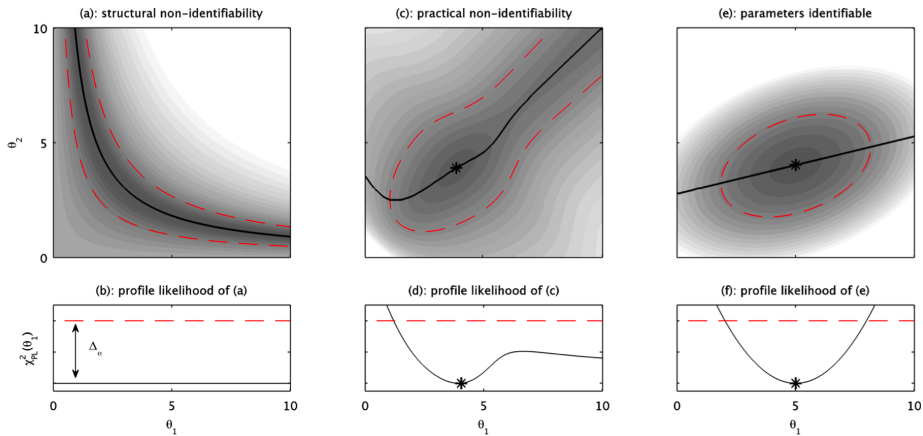


Figure: From Raue et al, 2010

T helper cell differentiation example

- ▶ Naive T cells sense the environment and react to different cytokine milieu and differentiate to functionally different T helper subsets
- ▶ Proper regulation and balance between various lymphocyte subsets is central the human immune system
- ▶ For example, so-called T helper 1 cells are involved in defense against intracellular bacteria
- ▶ Molecular mechanisms which control T helper cell differentiation are not at all fully understood
 - ▶ A more thorough understanding would allow e.g. better drug design to modulate immune response and help in autoimmune diseases

T helper cell differentiation example (2)

► Different T cell subsets

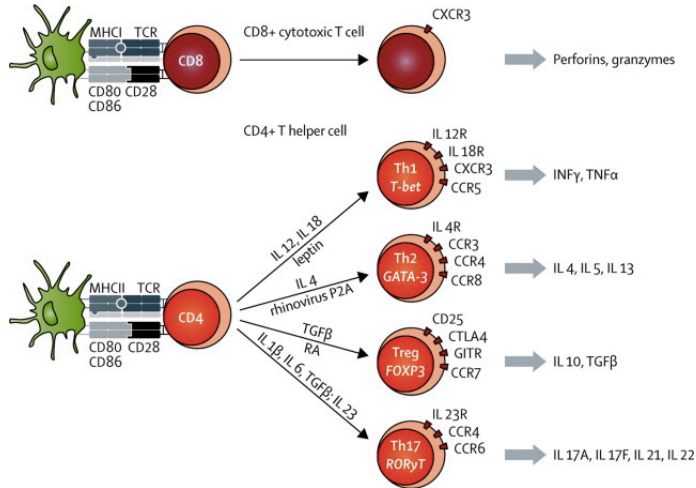


Figure: From (Brusselle et al, 2011)

T helper cell differentiation example (3)

- ▶ Key factors (i.e., genes) involved in and driving T helper 1 differentiation include
 - ▶ IFN- γ , T-bet, IL-12 and IL-12R β 2
- ▶ The standard/assumed molecular model, until recently, has been the following

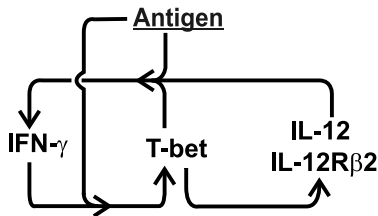


Figure: From (Schulz et al, 2009)

T helper cell differentiation example (4)

- The actual molecular mechanisms are assumed to follow the following ODE system

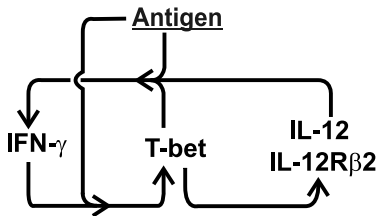


Figure: From (Schulz et al, 2009)

$$\frac{dTbet_{mRNA}}{dt} = \alpha_1 + f_{Tbet}(IFN_{Prot}, Rec_{Prot}, Ag) - \gamma_{Tbet} \cdot Tbet_{mRNA}$$

$$\frac{dTbet_{Prot}}{dt} = \beta \cdot Tbet_{mRNA} - \delta_{Tbet} \cdot Tbet_{Prot}$$

$$\frac{dRec_{mRNA}}{dt} = f_{Rec}(Tbet_{Prot}, Ag) - \gamma_{Rec} \cdot Rec_{mRNA}$$

$$\frac{dRec_{Prot}}{dt} = \beta \cdot Rec_{mRNA} - \delta_{Rec} \cdot Rec_{Prot}$$

$$\frac{dIFN_{mRNA}}{dt} = f_{IFN}(Tbet_{Prot}, Rec_{Prot}, Ag) - \gamma_{IFN} \cdot IFN_{mRNA}$$

$$\frac{dIFN_{Prot}}{dt} = \beta \cdot IFN_{mRNA} - \delta_{IFN} \cdot IFN_{Prot}$$

Figure: From (Schulz et al, 2009)

T helper cell differentiation example (5)

- ▶ To test the model, mRNA time-course measurements have been collected after inducing the Th1 differentiation
- ▶ ODE model is fit to the experimental data using simulated annealing

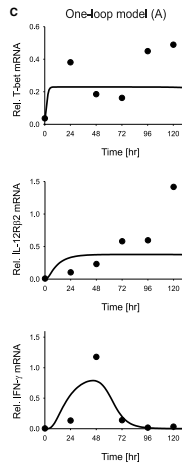
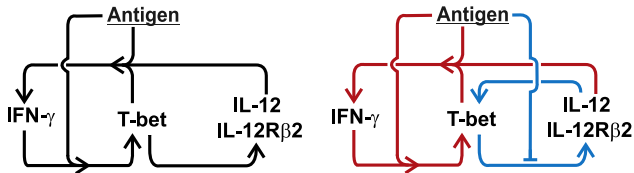


Figure: From (Schulz et al, 2009)

T helper cell differentiation example (6)

- Something wrong with the model? — another network model



$$f_{Tbet} = \alpha_2 \cdot \frac{Ag(t)}{K_1 + Ag(t)} \cdot \frac{IFN_{Prot}}{K_2 + IFN_{Prot}}$$

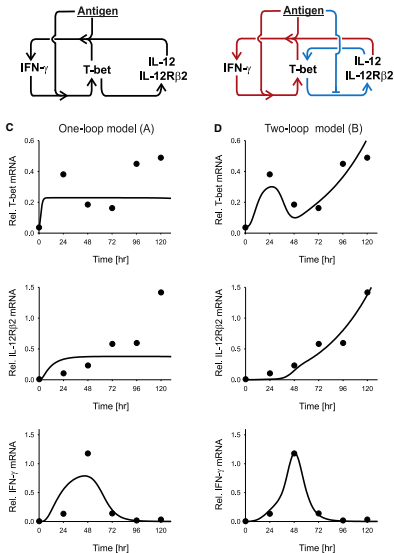
$$f_{Rec} = \alpha_4 \cdot \frac{Tbet_{Prot}}{K_8 + Tbet_{Prot}}$$

$$f_{Tbet} = \alpha_2 \cdot \frac{Ag(t)}{K_1 + Ag(t)} \cdot \frac{IFN_{Prot}}{K_2 + IFN_{Prot}} + \alpha_3 \cdot \frac{Rec_{Prot}}{K_3 + Rec_{Prot}}$$

$$f_{Rec} = \alpha_4 \cdot \frac{Tbet_{Prot}}{K_8 + Tbet_{Prot}} \cdot \frac{K_4}{K_4 + Ag(t)}$$

Figure: From (Schulz et al, 2009)

T helper cell differentiation example (7)



T helper cell differentiation example (8)

- The ODE model is predictive (at least qualitatively)

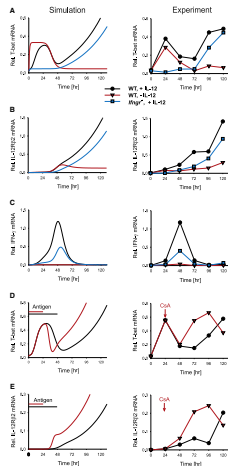


Figure: From (Schulz et al, 2009)

T helper cell differentiation example (9)

- More networks models — a fundamental question concerns choosing the correct biological network model

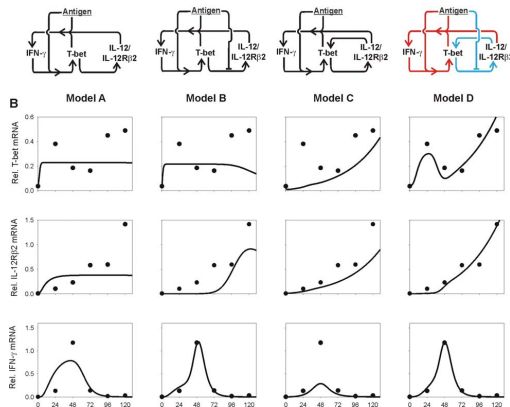


Figure: From (Schulz et al, 2009)

References

- ▶ Fröchlich et al. (2017) Scalable Inference of Ordinary Differential Equation Models of Biochemical Processes, *arXiv preprint arXiv:1711.08079*
- ▶ Raue et al. (2013) Lessons learned from quantitative dynamical modeling in systems biology, *PLoS One*, 8(9), e74335.
- ▶ Raue et al. (2010) Identifiability and observability analysis for experimental design in nonlinear dynamical models, *Chaos*, 20, 045105.
- ▶ Brusselle GG et al, (2011) New insights into the immunology of chronic obstructive pulmonary disease, *Lancet*, 378(9795):10–16.
- ▶ Schulz et al, (2009) Sequential Polarization and Imprinting of Type 1 T Helper Lymphocytes by Interferon-g and Interleukin-12, *Immunity*, 30, 673–683.