# Inferring rate coefficients of biochemical reactions from noisy data with KInfer

Paola Lecca

*The Microsoft Research - University of Trento*
*Centre for Computational and Systems Biology*

lecca@cosbi.eu

Alida Palmisano

*The Microsoft Research - University of Trento*
*Centre for Computational and Systems Biology*
*DISI - University of Trento*

palmisano@cosbi.eu

Corrado Priami

*The Microsoft Research - University of Trento*
*Centre for Computational and Systems Biology*
*DISI - University of Trento*

priami@cosbi.eu

# Inferring rate coefficents of biochemical reactions from noisy data with KInfer

Paola Lecca, Alida Palmisano and Corrado Priami

**Abstract**

Dynamical models of inter- and intra-cellular processes contain the rate constants of the biochemical reactions. These kinetic parameters are often not accessible directly through experiments, but they can be inferred from time-resolved data. Time resolved data, that is, measurements of reactant concentration at series of time points, are usually affected by different types of error, whose source can be both experimental and biological. The noise in the input data makes the estimation of the model parameters a very difficult task, as if the inference method is not sufficiently robust to the noise, the resulting estimates are not reliable. Therefore "noise-robust" methods that estimate rate constants with the maximum precision and accuracy are needed. In this report we present the probabilistic generative model of parameter inference implemented by the software prototype KInfer and we show the ability of this tool of estimating the rate coefficients of models of biochemical network with a good accuracy even from very noisy input data.

## 1  Introduction

The relation between the instantaneous rate of reaction and the concentrations of the reactants at any moment is given by the law of mass action: i.e. the rate at which a substance takes part in a reaction is proportional to its concentration raised to a power which represents the number of molecules taking part in the reaction. The ability to infer these constants of proportionality for a system of biochemical reactions is crucial in systems biology, yet their direct measurement is a challenging experimental problem.

Parameter estimation is commonly achieved by the best fit of numerical simulations to experimental observations. The fitting procedure is based on optimization techniques where a measure of the distance between model prediction and experimental data (the cost function) is used as the optimality criterion to be minimized. In most approaches dealing with parameter estimation the cost function is the likelihood function, also know as joint transitional density. It expresses the probability of obtaining the observed outcomes in terms of measured systems variables and parameters. Thus it can be used to determine unknown parameters based on known outcomes. The optimal values of the parameters can be estimated by maximizing the likelihood function (maximum likelihood criterion) or, equivalently, by minimizing the log-likelihood function. However, when estimating parameters of dynamical systems with optimization methods a number of difficulties may arise, the main of which are convergence to local solutions, very flat objective function in the neighborhood of the solution, over-determined models, and non-differentiable terms in the systems dynamics. Due to the non-linear nature of the dynamics of the biological processes, these problems are often multimodal, so that traditional gradient based methods fail to identify the global solution and may converge to a local minimum. Moreover, in the case in which a bad fit has been performed, there is no way of knowing if it is due to a wrong model structure or if it is consequent to a local convergence.

The recent literature reports many examples of new effective methods attempting either to work out these difficulties or to develop new methodologies of parameter estimation both in deterministic and stochastic models. Here we briefly mention the most recent ones. Polisetty et al. in [12] suggested global optimization techniques as alternative to traditional local methods. Rodrigez-Fernandez et al. in [14] developed a hybrid stochastic-deterministic global optimization method. Moles et al. in [11] explored several state-of-the-art

deterministic and stochastic global optimization techniques and compared their accuracy and effectiveness on nonlinear biochemical dynamic models. Tian et al. [4] presented simulated maximum likelihood method to evaluate parameters in stochastic models described by stochastic differential equation. They propose different types of transitional probability and a genetic optimization algorithm to search for optimal reaction rates. Chou et al. [2] developed an alternate regression method, that dissects the parameter inference problem into iterative steps of linear regression. Sugimoto et al. [15] provided a computational technique based on genetic programming that simultaneously generates biochemical equations and their parameters from time series data. Reinker et al. [13] are the authors of the approximate maximum likelihood method and the singular value decomposition likelihood method that estimate stochastic reaction constants from molecule count data measured with errors at discrete time points. Tools for parameter fitting through regression or maximum likelihood methods can be found as integral part of simulation tools (e. g. Copasi [8]), but there exist also stand-alone softwares exclusively designed for that purpose, like PET [19]. Finally, we mention the works of Boys [1], Golitki [5] and Wilkinson [17], that developed Bayesian model-based inference techniques specific for discrete models. Bayesian scheme depart from the approaches previously mentioned. They offer some advantages over the maximum likelihood methods, for instance when the volume of data is limited or the analytic form of the kinetic model makes the maximization of the likelihood not straightforward. The disadvantages of the most part of the current tools for parameter estimation is the lack of robustness to the noise and the absence of any estimate of experimental error in their outcome. Experimental uncertainties on parameters propagate from the measurements of the concentrations of the species. Returning the parameters with an estimate of their uncertainty is essential if we want to use the tool in the context of optimal experimental design. Moreover, the most part of the current tools, based on optimization techniques suffer from the problem of univocally finding the solution global optimization, and ask the user to provide *a priori* the optimization algorithm with the region of parameter space in which to perform the search for the global max/minimum.

In this paper we present a novel approach to the model calibration, that proposes the solution to these difficulties and whose estimation accuracy is robust w. r. t. the experimental noise. The method is based on a probabilistic, generative model of the variations in reactant concentration. Given reactant species, we observe time series concentrations for each of the species, gathered in $N$ state vectors $\mathbf{X}_1, \ldots, \mathbf{X}_N$ , our method discretizes the law of mass action and provides a tool to predict the values of the variables $\mathbf{X}_i$ at time $t$ , conditioned on their values at the previous time point. The variations of the concentration of the species at different time points are conditionally independent by the Markov nature of the discrete model of the law of mass action. Assuming the observation noise to be Gaussian with variance $\sigma^2$, the probability of observing a variation $D_i$ for the concentration $X_i$ of species $i$ between time $t_{k-1}$ and $t_k$ is a Gaussian with variance depending on $\sigma$ and mean the expectation value of the law of mass action under the noise distribution. The discretization of the law of mass action provides a model for the variations of the species concentration, rather than a model for the time-trajectory of the species concentrations. This makes the evaluation of the expectation value of law mass action function (the integral of the transitional probability) simpler and analytically tractable. The rate coefficients and the level of noise are then obtained by maximizing the likelihood function defined by the observed variations. Our method returns the rate coefficients, the level of noise $\sigma$ and an error range on the estimates of rate constants. Its probabilistic formulation is key to a principled handling of the noise inherent in biological data, and it allows for a number of further extensions, such as a fully Bayesian treatment of the parameter inference and automated model selection strategies based on the comparison between marginal likelihoods of different models. Finally, the implementation of this method may be used as an interface tool, connecting the outcomes of the wet-lab activity for the concentration measurements and the software for the simulation of chemical kinetics.

We show the ability of our algorithm of obtaining reasonable estimates for the rate coefficients in case studies of different complexity. In particular we present the results of the application of KInfer, the software that we developed and that implements the our inference procedure (on synthetic and real data) to the following case studies: gene transcription and expression, gene transcriptional regulation, and coupled autocatalytic reactions of Lotka model, and spike generation in neuronal dynamics. The parameters of the kinetics of these pathways are known, since they were experimentally determined and widely documented

in literature. For the last case study, in particular, we possess real experimental time course data. Thus, we could compare our estimates of the parameters with the known values to assess the soundness of the methodology and the performance of its implementation. This work is intended to be a *statement* paper of our inference procedure and of its usefulness in an experimental context. We do not report here the comparative analysis of the accuracy of our method with respect to the others, but for each case study we indicate the literature reference in which the reader can find the estimates obtained by other methods and we recall that KInfer is downloadable for free at http://www.cosbi.eu.

## 2  The model

Consider $N$ reactant species, $S_1, S_2, \ldots, S_N$, with concentrations $X_1, X_2, \ldots, X_N$, that evolve according to a system of rate equations

$$\frac{dX_i}{dt} = f_i(\mathbf{X}^{(i)}(t); \theta_i) \tag{1}$$

where $\theta_i$ , $i = 1, 2, \ldots, N$ , is the vector of the rate coefficients, which are present in the expression of the function $f_i$ . We wish to estimate the set of parameters $\mathbf{\Theta} = \cup \theta_i$ $(i = 1, 2, \ldots, N)$, whose element $\theta_i$ is the set of rate coefficients appearing in the rate equations of i-th species, therefore

$$\theta_1 = \{\theta_{11}, \theta_{12}, \ldots, \theta_{1N_1}\}, \ldots, \theta_N = \{\theta_{N1}, \theta_{N2}, \ldots, \theta_{NN_N}\}$$

$\mathbf{X}^{(i)}$ is the vector of concentrations of chemicals that are present in the expression of the function $f_i$ for the species $i$. According to the law of mass action, the functions $f_i$ have the general form

$$f_i(\mathbf{X}^{(i)}(t); \theta_i) =$$

$$= \theta_{i1} \prod_{w \in S_1 \subseteq [1,N]} X_w^{\alpha_w} + \cdots + \theta_{iN_i} \prod_{w \in S_{N_i} \subseteq [1,N]} X_w^{\alpha_w} = \sum_{h=1}^{N_i} \left( \theta_{ih} \prod_{w \in S_h} X_w^{\alpha_w} \right) \tag{2}$$

where $\alpha_w \in \mathbf{R}$, and $N_i$ is the number of parameter in the $f_i$ rate equation The rate equations in (2) form the so-called Generalized Mass Action law. We assume we have noisy observations $\hat{X}_i = X_i + \epsilon$ at times $t_0, \ldots, t_M$ , where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is a Gaussian noise term with mean zero and variance $\sigma$. With this choice we are assuming that the concentration measurements are not significantly affected by systematic errors, but by uncontrolled random errors and that an error is equally likely to occur in either positive or negative direction with respect to the symmetry axis of the distribution.

We also assume a number $M$ of concentration measurements for each considered species. Approximating the rate equation (1) as a finite difference equation between the observation times, gives

$$X_i(t_k) = X_i(t_{k-1}) + (t_k - t_{k-1}) f_i(\mathbf{X}^{(i)}(t_{k-1}); \theta_i) \tag{3}$$

where $k = 1, \ldots, M$. In Eq. (3) the rate equation is viewed as a model of increments/decrements of reactant concentrations; i.e., given a value of the variables at time $t_{k-1}$ , the model can be used to predict the value at the next time point $t_k$. Increments/decrements between different time points are conditionally independent by the Markov nature of the model (3). Therefore, given the Gaussian model for the noise, it is possible to estimate the probability to observe the value $\hat{X}_i(t_k)$ given the model at time $t_{k-1}$, $X_i(t_{k-1})$, and the set of parameters $\theta_i$, as

$$p\left(\hat{X}_i(t_{k-1}) | X_i(t_{k-1})\right) = \mathcal{N}\left(X_i(t_{k-1}) + (t_k - t_{k-1}) f_i(X_i(t_{k-1}, \theta_i)), \sigma^2\right) \tag{4}$$

We then also have that the true value of $X_i(t_k)$ is normally distributed around the observed value $\hat{X}_i(t_k)$, so that

$$p\Big(X_i(t_{k-1})|\hat{X}_i(t_{k-1}))\Big) = \mathcal{N}\Big(\hat{X}_i(t_{k-1}), \sigma^2\Big) = \tag{5}$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\Big[ -\frac{(X_i(t_{k-1}) - \hat{X}_i(t_{k-1}))^2}{2\sigma^2} \Big]$$

Therefore, the probability to observe a variation $D_i(t_k) = X_i(t_k) - X_i(t_{k-1})$ for the concentration of the $i$-th species between the time $t_{k-1}$ and $t_k$, given the parameter vector $\theta_i$ is

$$p(D_i(t_k)|\theta_i, \sigma) = \mathcal{N}\Big(E\big[f_i(\mathbf{X}^{(i)}(t_{k-1}), \theta_i)\big], 2\sigma^2\Big) \tag{6}$$

and

$$E\big[f_i(\mathbf{X}^{(i)}(t_{k-1}, \theta_i))\big] = \int_{\Omega_{\mathbf{X}^{(i)}}} f_i(\mathbf{X}^{(i)}(t_{k-1}), \theta_i) \prod_{i=1}^{K_i} \Big[p_i\Big(X_i(t_{k-1})|\hat{X}_i(t_{k-1})\Big)\Big] d\mathbf{X}^{(i)} \tag{7}$$

where $\Omega_{\mathbf{X}^{(i)}}$ is the sample space of $\mathbf{X}^{(i)}$, and $K_i$ is the number of chemical species in the expression for $f_i$. While the increments/decrements are conditionally independent given the starting point $X_i(t_k)$, the random variables $D_i(t_k)$ are not independent of each other. Intuitively, if $X_i(t_k)$ happens to be below its expected value because of random fluctuations, then the following increment $D_i(t_{k+1})$ can be expected to be bigger as a result, while the previous one $D_i(t_k)$ will be smaller. A simple calculation allows us to obtain the covariance matrix of the vector of increments for the $i$-th species. This is a banded matrix $\mathbf{C}_i \equiv \mathbf{C} = \text{Cov}(\mathbf{D}_i)$ with diagonal elements given by

$$E\Big[D_i^2(t_k) - E[D_i^2(t_k)]\Big] = 2\sigma^2$$

and a non-zero band above and below the diagonal given by

$$E\Big[\big(D_i(t_k) - E[D_i(t_k)]\big)\big(D_i(t_{k-1}) - E[D_i(t_{k-1})]\big)\Big] = -\sigma^2$$

with all other entries zero. The likelihood for the observed increments/decrements therefore will be

$$p(\mathbf{D}|\mathbf{\Theta}) = \prod_{i=1}^{N} \mathcal{N}(\mathbf{D}_i|\mathbf{m}_i(\mathbf{\Theta}), \mathbf{C}) = \Big(\frac{1}{\sqrt{2\pi \det(\mathbf{C})}}\Big)^N e^{\sum_{i=1}^{N} -\frac{1}{2}(\mathbf{D}_i - \mathbf{m}_i)^T \mathbf{C}^{-1}(\mathbf{D}_i - \mathbf{m}_i)} \tag{8}$$

where $\mathbf{D} = \{\mathbf{D}_1, \ldots, \mathbf{D}_N\}$, $\mathbf{D}_i = D_i(t_1), D_i(t_2), \ldots D_i(t_M)$ $(i = 1, 2, \ldots, N)$, and $\mathbf{m}_i(t_{k-1}) \equiv E\Big[f_i(\mathbf{X}(t_{k-1}), \theta_i)\Big]$.

The Eq. (8) can be optimized w. r. t. the parameters $\mathbf{\Theta} = (\theta_1, \theta_2, \ldots, \theta_N)$ of the model to yield estimates of the parameters themselves and of the noise level. The chief numerical problem of this approach is the computation of the expectations of the rate functions given by equation (7). Non-integer values of the coefficients $\alpha$ can make estimating the integral analytically difficult. We propose an approximate method in which the Gaussian noise is replaced by an approximate uniform (white) noise, with the amplitude of the uniform noise being obtained as a sample from the Gaussian cumulative distribution function. At the first order, for small $\sigma$, we can approximate the Gaussian with zero mean and variance $\sigma$ with an uniform distribution defined on the interval $[-\frac{\sqrt{2\pi}\sigma}{4}, \frac{\sqrt{2\pi}\sigma}{4}]$, so that

$$\prod_{i=1}^{K_i} p_i = \prod_{i=1}^{K_i} \chi_i \tag{9}$$

where

$$\chi_i(X_i) = \begin{cases} \frac{2}{\sqrt{2\pi}\sigma} & \text{if } -\frac{\sqrt{2\pi}\sigma}{4} \leq X_i \leq \frac{\sqrt{2\pi}\sigma}{4} \\ 0 & \text{otherwise.} \end{cases}$$

4

This approximation makes the calculation of the expectation value of the rate equation (Eq. (7)) simpler and reduces the computational time of the procedure. Moreover, experiments not illustrated in this paper demonstrate that it does not influence the accuracy of the parameter estimates until $\sigma$ is less that 30% of the concentration measurement.

Substituting Eq. (9) in Eq, (7) gives

$$E[f_i(\mathbf{X}^{(i)}(t_{k-1}),\theta)] = \left(\frac{2}{\sqrt{2\pi}\sigma}\right)^{K_i} \int_{\hat{X}-\frac{\sqrt{2\pi}\sigma}{4}}^{\hat{X}+\frac{\sqrt{2\pi}\sigma}{4}} f_i(\mathbf{X}^{(i)}(t_{k-1}),\theta_i)d\mathbf{X}^{(i)} \tag{10}$$

Now, substituting Eq. (2) in Eq. (10) leads to

$$E[f_i(\mathbf{X}^{(i)}(t_{k-1}),\theta_i)] =$$
$$= \left(\frac{2}{\sqrt{2\pi}\sigma}\right)^{K_i} \left\{ \sum_{h=1}^{N_i} \theta_{ih} \left[ \left(\frac{\sqrt{2\pi}\sigma}{2}\right)^{\#(S-S_h)} \times \right. \right.$$
$$\left. \left. \times \prod_{w \in S_h} \frac{1}{\alpha_w+1} \left( \left(\hat{X}_w + \frac{\sqrt{2\pi}\sigma}{4}\right)^{\alpha_w+1} - \left(\hat{X}_w - \frac{\sqrt{2\pi}\sigma}{4}\right)^{\alpha_w+1} \right) \right] \right\} \tag{11}$$

where $S$ is the set containing the indexes referring to all the $K_i$ species appearing in $f_i$, and $\alpha_w \neq -1$. In case some orders are equal to -1 Eq. (11) takes the following form

$$E[f_i(\mathbf{X}^{(i)}(t_{k-1}),\theta_i)] = \left(\frac{2}{\sqrt{2\pi}\sigma}\right)^{K_i} \sum_{h=1}^{N_i} \theta_{ih} \left\{ \left(\frac{\sqrt{2\pi}\sigma}{2}\right)^{\#(S-S_h)} \times \right.$$
$$\times \left[ \prod_{w \in S'_h} \frac{1}{\alpha_w+1} \left( \left(\hat{X}_w + \frac{\sqrt{2\pi}\sigma}{4}\right)^{\alpha_w+1} - \left(\hat{X}_w - \frac{\sqrt{2\pi}\sigma}{4}\right)^{\alpha_w+1} \right) \right] \times$$
$$\times \left. \left[ \prod_{w \in S''_h} \ln \frac{\hat{X}_w + \frac{\sqrt{2\pi}\sigma}{4}}{\hat{X}_w - \frac{\sqrt{2\pi}\sigma}{4}} \right] \right\} \tag{12}$$

where $S'_h$ is the set of indexes $\{h'_1, h'_2, \ldots, h'_s\}$ such that $\alpha_{h'} \neq -1 \ \forall h' \in S'_h$, and $S''_h$ is the set of indexes $\{h''_1, h''_2, \ldots, h''_s\}$ such that $\alpha_{h''} = -1 \ \forall h'' \in S''_h$.

If in the Eq. (8), $\mathbf{m}_i$ is substituted with the expression (11) or (12), Eq. (8) becomes more tractable and can be optimized w. r. t. the parameters $\mathbf{\Theta} = (\theta_1, \theta_2, \ldots, \theta_N)$ and $\sigma$. The values of the model's parameters for which $p(\mathbf{D}|\mathbf{\Theta})$ has a maximum are the most likely values giving the observed kinetics.

## 3 KInfer: a prototype for parameter inference

We developed the prototype KInfer (Kinetics Inference), that implements the procedure described in the previous section. The tool consists of four main blocks: 1) the input interface, 2) the model generator, 3) the maximization algorithm and 4) the output interface (Fig. 1). A view of the screenshots of the tool is shown in Fig. 2.

The specification of the reactions must end with semicolon. Along with the specification of the set of reactions involved in the system, KInfer requires the experimental time series data, in tabular text format, of the concentration (or number of molecules) of the species present in the system. The option "Load concentrations" in the File menu of the front-end allows the user to download the experimental times series of concentrations. From the set of chemical reactions the tool automatically generates the ordinary differential equations model, consisting of a system of equations of the form of Eq. (2) (see the field "Automatic model"
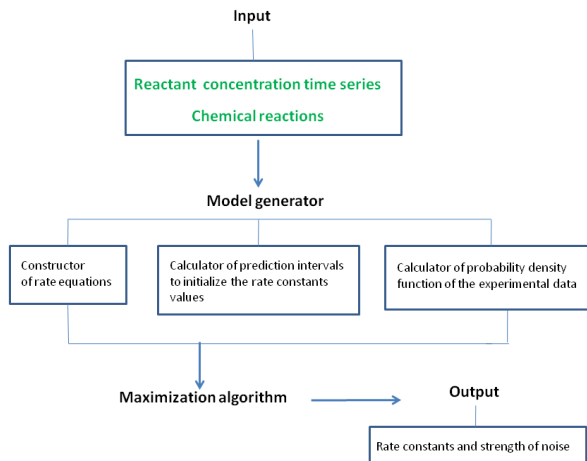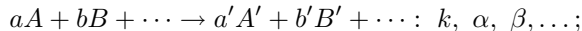
Figure 1: main modules of KInfer.

The tool takes as input the set of chemical reactions describing the kinetics of the systems, specified in the following syntax

$$aA + bB + \cdots \rightarrow a'A' + b'B' + \cdots : \ k, \ \alpha, \ \beta, \ldots ;$$

On the left-hand side of the arrow, the reactants $(A, B, \ldots )$ and the reactants stoichiometric coefficients $(a, b, \ldots )$ are indicated, whereas on the right-hand side the products $(A', B', \ldots )$ and the product stoichiometric coefficients $(a', b', \ldots )$ are indicated. The reaction specification contains also the indication of the name of the rate constant after colon and the partial orders of reaction $(\alpha, \beta, \ldots )$.
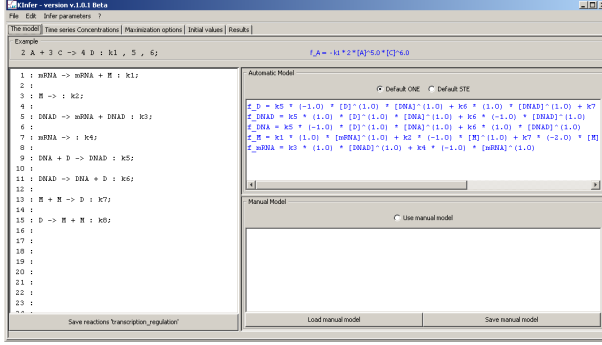
in the front-end in Fig. 2). However, the user is allowed to insert a different model that can be entered in the "Manual Model" part of the interface. The user is allowed also to enter an ordinary differential equation model without specifying the reaction in the standard "chemical" notation. The tool processes the inputs and it derives from the data set of the concentration time-series and from the model of rate equation the form of the probability density function in Eq. (8) to maximize and the initial guesses for the parameters. Although the tool automatically calculates the initial guesses of the parameters, the user is allowed to change the estimated values as well as to directly insert new different estimates.

The optimization algorithm of KInfer is a Genetic Algorithm [3]. This choice has been driven by the fact that a biological model of realistic size and complexity presents a high number of parameters with possible nonlinear relations between them. For technical details we refer the reader to [3]. Here we simply recall that a genetic algorithm is a population based stochastic optimization technique, that, starting from a set of initial guesses about the solution, determines the next set of possible solutions to the optimization problem on the basis of the results obtained from the preceding set and approaching step by step, as in a typical evolutionary scenario, toward the better solution. These methods have been designed primarily to address problems that cannot be tackled through traditional optimization algorithms. Such problems are characterized by discontinuities, lack of derivative information, noisy function values and disjoint search spaces.
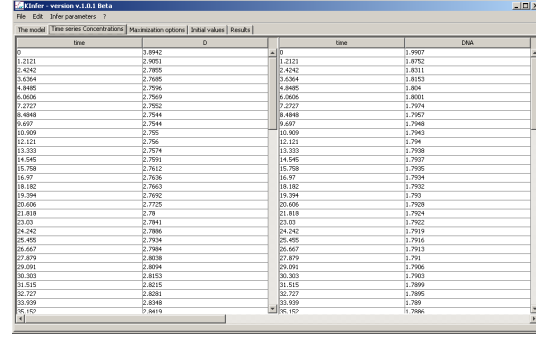
The search for the optimal values of rate constants can be made more efficient if we provide the algorithm of optimization of Eq. (8) with the initial guesses for these constants. In this way the algorithm does not waste time in exploring large regions of the parameter space or regions in which the model in Eq. (3) is not valid. For this purpose, we also developed and included in KInfer a procedure for the automatic calculation of the initial guesses of the parameters. Therefore, the task to direct the inference method to efficiently exploring the parameter space is not left to the user, that often does not have a precise idea about a reasonable value of the parameters. Because of lack of space, we refer the reader to [10], where the entire procedure is explained, adn we report here the basic ideas. The derivatives $dX/dt$ at all measured time points $t_k$ can be interpreted as slopes. Given the species $i$ (with $i = 1, \ldots, N$), we can estimate these slopes from the data as $s_i(t_k)$, and approximate the differential equations as

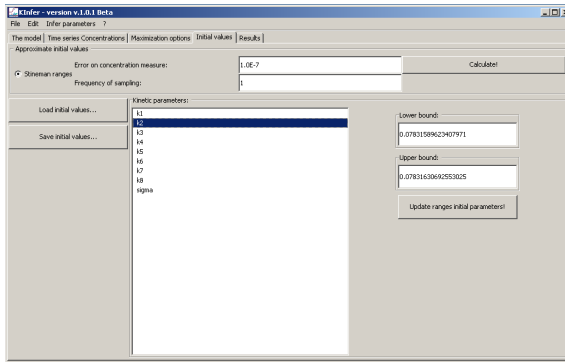$$s_i(t_k) \approx \left. \frac{dX_i}{dt} \right|_{t=t_k} \tag{13}$$

If the data consist of $N$ species and the concentration of each species $i$ is measured at $M$ time points
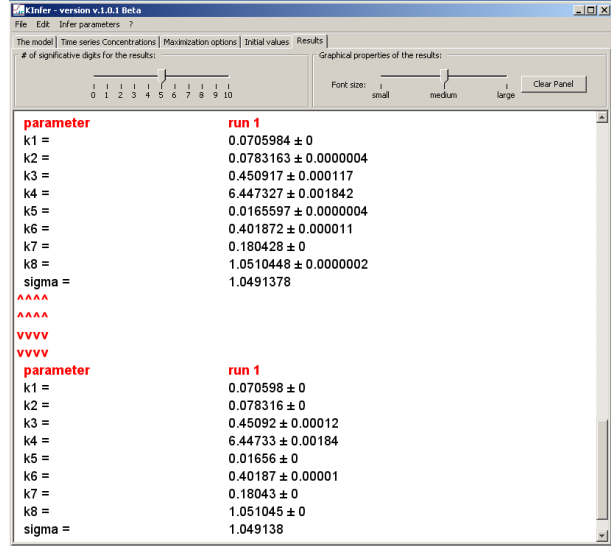
6

Figure 2: (a) Front-end of KInfer. It is divided in three regions: let us call them "region A" on the left, "region B" and 'region C" on the right. In region A the user can write the reactions of the system. In region B the rate equation model is automatically generated. If the user wish to change these equations, he is allowed to do it in region C, where he can write a new different set of rate equations. (b) The table of uploaded experimental data. (d) The settings of the estimator of the parameter initial guesses, and (c) the screenshot showing the results.

$(X(t_1), X_i(t_2), \ldots, X_i(t_M))$, we estimate $M \times N$ slopes $s_i(t_k)$ $(k = 1, \ldots, M)$. In fact, for each species we have $M$ differential equation of the form

$$s_i(t_k) \approx f_i(X_1(t_k), X_2(t_k), \ldots, X_N(t_k); \theta_{i1}, \theta_{i2}, \ldots, \theta_{iN_i}) \tag{14}$$

that form a system of $M$ algebraic equations with $M \times N_i$ unknown variables $\theta$s, as the slopes $s$ are measurable from the data. In general, $M \gg N_i$, so that the system of $M \times N_i$ equation results overdetermined. In order to avoid this situation without having recourse to techniques of least squares fitting for overdetermined systems, we re-sample the experimental time course of the species $i$ order to have $M = N_i$ and a good approximation of the original curve interpolating all the $M$ points. At this stage we are not interested in a very precise estimate of the rate constants, but only in an approximate guess. A system of equations similar to the system (15) can be written also for the experimental uncertainties $\Delta s_i$ affecting the slopes $s_i$:

7

$$\Delta s_i(t_k) \approx \Delta f_i(X_1(t_k), X_2(t_k), \ldots, X_N(t_k); \theta_{i1}, \theta_{i2}, \ldots, \theta_{iN_i}) \tag{15}$$

where

$$\Delta s_i = \Delta\left[\theta_{i1} \prod_{j \in S_1 \subseteq [1,N]} X_j^{\alpha_j}\right] + \Delta\left[\theta_{i2} \prod_{j \in S_2 \subseteq [1,N]} X_j^{\alpha_j}\right] + \cdots + \Delta\left[\theta_{iN_i} \prod_{j \in S_{N_i} \subseteq [1,N]} X_j^{\alpha_j}\right] \tag{16}$$

By using the standard formulas of the error propagation, a single term of the sum on the right-hand side of Eq. (16) is found to be

$$\frac{\Delta\left[\theta_{i1} \prod_{j \in S_1 \subseteq [1,N]} X_j^{\alpha_j}\right]}{\left|\theta_{i1} \prod_{j \in S_1 \subseteq [1,N]} X_j^{\alpha_j}\right|} = \frac{\Delta\theta_{i1}}{\theta_{i1}} + \frac{\Delta\left[\prod_{j \in S_1 \subseteq [1,N]} X_j^{\alpha_j}\right]}{\left|\prod_{j \in S_1 \subseteq [1,N]} X_j^{\alpha_j}\right|} = \frac{\Delta\theta_{i1}}{\theta_{i1}} + \sum_{h=1}^{\#S_1} |\alpha_h| \frac{\Delta X_h}{|X_h|}$$

where $\#S_1$ is the cardinality of the set $S_1$. Therefore, Eq. (16) becomes

$$\Delta s_i = \sum_{\nu=1}^{N_i} \left\{ \left( \frac{\Delta\theta_{i\nu}}{\theta_{i\mu}} + \sum_{h=1}^{\#S_\nu} |\alpha_h| \frac{\Delta X_h}{|X_h|} \right) \cdot \left| \theta_{i\mu} \prod_{j \in S_\nu \subseteq [1,N]} X_j^{\alpha_j} \right| \right\} \tag{17}$$

Now, assuming that the measurements of times are not affected by errors, the error $\Delta s_i$ is calculated from Eq. (3) as follows

$$\Delta s_i(t_k) = \frac{1}{t_k - t_{k-1}} \Big( \Delta X_i(t_k) - \Delta X_i(t_{k-1}) \Big)$$

where $\Delta X_i(t_k)$ is the experimental error on the measurement of concentration of species $i$ at time $t_k$. Therefore $\Delta s_i(t_k)$ can be obtained from the data, and the system (17) can be solved, with the same procedure used for the system (15), to find the size of the prediction intervals of the $\theta$s: the $\Delta\theta$. These intervals are also approximate measures of the errors that from the concentration measurements propagate to the rate constants.

## 4 Case studies

Here we provide some validation tests on biochemical networks typically considered in the literature concerning parameter inference. For each case study we provide a picture of the reaction network, the table comparing actual and estimated parameters, and the simulation curves of the system dynamics obtained with the actual and estimated parameters. We did not include in this manuscript the experimental and/or synthetic time series of the concentrations we used as input of our procedure to infer the parameter. For shortness reasons, we simply report the time resolution and the number of data points. The errors on the parameter estimates computed by our procedure are neither computational errors imputable to the precision of the integration and optimization algorithms or to the variance of repeated inference runs. They are experimental errors that propagate from the concentration measurements to the model rate coefficients. Therefore, their values are not comparable with the values of errors on the parameters estimates obtained by the references cited in each case study, which we refer the reader to. The results that we obtained confirm that the procedure converges to the expected solution within the experimental errors and the strength of noise affecting the input data. Some discrepancies between the actual value and its estimate in the reported case studies are mainly due to the level of noise and to the approximations introduced by the discretization of the rate equations. Nevertheless, in most cases these discrepancies are found in parameters to which the model is not sensitive, and thus, the its dynamics is not strongly influenced by them.

### 4.1 Case study 1: a didactic example of biochemical network

The system depicted in Fig. 3 is representative of a small biochemical network of 4 interacting species. The network has two feedback loops: 1. the species $X_5$ inhibits the production of species $X_1$, and 2. the

species $X_4$ promotes the activation of $X_5$. A numerical implementation with typical parameters is given by the set of ordinary differential equations in Fig. 3. This system of equations has been used to create the artificial time series of 51 data points with a time resolution of 0.2. Typical units might mM for the concentration and minutes for the times, but the example could as well run on an hourly scale and with variables of different nature. Table 1 lists the results and Fig. 4 shows the dynamic simulations. Within the experimental uncertainties, these results are in agreement with the expected ones and with those in [2].



$$X_1' = \theta_1 X_3^{-0.8} - \theta_2 X_1^{0.5} \qquad X_1(t_0) = 1.4$$
$$X_2' = \theta_3 X_1^{0.5} - \theta_4 X_2^{0.75} \qquad X_2(t_0) = 2.7$$
$$X_3' = \theta_5 X_1^{0.75} - \theta_6 X_3^{0.5} X_4^{0.2} \qquad X_3(t_0) = 1.2$$
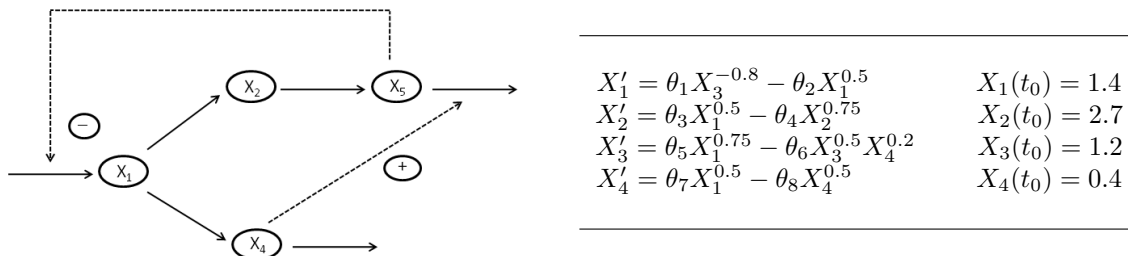$$X_4' = \theta_7 X_1^{0.5} - \theta_8 X_4^{0.5} \qquad X_4(t_0) = 0.4$$

Figure 3: a didactic example of biochemical network with four variables and the system of ordinary differential equations describing it. The concentration of $X_5$ is supposed to be constant.

| Parameter | Actual value | Initial guesses | Estimated value |
|---|---|---|---|
| $\theta_1$ | 12 | [10.18; 13.84] | $11.37 \pm 3.66$ |
| $\theta_2$ | 10 | [8.28; 11.74] | $9.39 \pm 3.46$ |
| $\theta_3$ | 8 | [9.81; 9.87] | $9.83 \pm 0.06$ |
| $\theta_4$ | 3 | [3.92; 3.99] | $3.98 \pm 0.07$ |
| $\theta_5$ | 3 | [2.91; 2.96] | $2.94 \pm 0.05$ |
| $\theta_6$ | 5 | [4.89; 4.91] | $4.90 \pm 0.02$ |
| $\theta_7$ | 2 | [1.50; 2.55] | $1.84 \pm 1.05$ |
| $\theta_8$ | 6 | [4.01; 8.17] | $5.5 \pm 4.16$ |
| $\sigma$ | 0.1 | [0.1; 0.3] | 0.3 |

Table 1: Case study 1: estimated parameter values for the network in Fig. 3

## 4.2   Case study 2: gene transcription and transcriptional regulation

In this test, we first consider the transcription of a single gene as given by the model of Golding et al. in [7, 14]. The DNA for the tagged mRNA is switched on and off by polymerase binding and unbinding, respectively. Only polymerase-bound DNA is transcribed into mRNA. The system is depicted in Fig. 5. We set the initial conditions $DNA_{OFF} = 1$, $DNA_{ON} = 0$ and $mRNA = 0$, and we generated a set of 100 data points at at temporal resolution of 1. Typical measurements units are "number of molecules" for the amount of the species and "minutes" for the time. Our estimates of parameters are reported in Table 2; the comparison of the estimated and experimental system's behavior in Fig. 6 shows a strong agreement. The accuracy of the results is comparable with one of those obtained by Reinker et al. [13] for the same network.

| Parameter | Actual value | Initial guesses | Estimated value |
|---|---|---|---|
| $\theta_1$ | 0.027 | [0.0242; 0.0249] | $0.0244 \pm 0.0007$ |
| $\theta_2$ | 0.1667 | [0.151; 0.152] | $0.152 \pm 0.001$ |
| $\theta_3$ | 0.4 | [1.578; 2.385] | $1.579 \pm 0.807$ |
| $\sigma$ | 0.5 | [0; 1] | 0.445 |

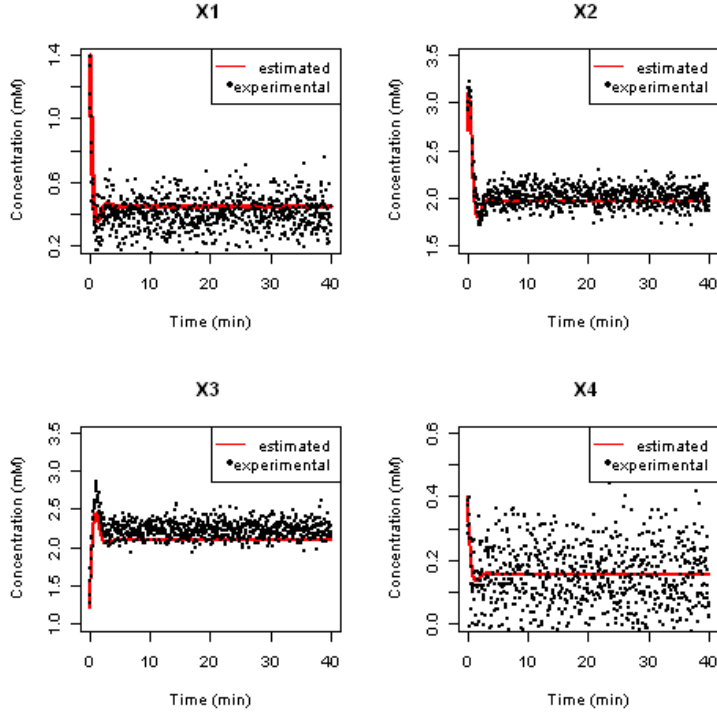Table 2: Case study 2: estimated parameter values for the Golding's model of gene transcription (Fig. 5).

Figure 4: simulations of case study 1. Experimental and estimated time behavior of the species of the biochemical network in Fig. 3.



$$DNA_{OFF} \xrightarrow{\theta_1} DNA_{ON}$$
$$DNA_{ON} \xrightarrow{\theta_2} DNA_{OFF}$$
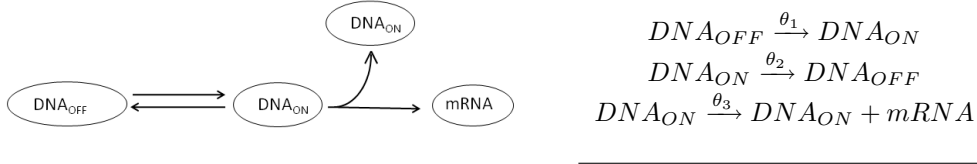$$DNA_{ON} \xrightarrow{\theta_3} DNA_{ON} + mRNA$$

Figure 5: Golding's model of gene transcription process.

Table 3 reports the estimates of the rate constants, that within the estimated error ranges, are in agreement with the actual values and with the results obtained by [6]. Figure 8 shows the actual and the estimated dynamics.

| Parameter | Actual value | Initial guesses | Estimated value |
|-----------|-------------|-----------------|-----------------|
| $\theta_1$ | 0.043 | [0.01; 0.08] | $0.042 \pm 0.007$ |
| $\theta_2$ | 0.0007 | [0.0001; 0.001] | $0.0004 \pm 0.0009$ |
| $\theta_3$ | 0.715 | [0; 1] | $0.1051 \pm 0.1$ |
| $\theta_4$ | 0.00395 | [0.00340; 0.00386] | $0.0038 \pm 0.0005$ |
| $\theta_5$ | 0.02 | [0.01; 0.04] | $0.019 \pm 0.03$ |
| $\theta_6$ | 0.4791 | [0; 1] | $0.62 \pm 0.17$ |
| $\theta_7$ | 0.083 | [0.01; 0.2] | $0.12 \pm 0.02$ |
| $\theta_8$ | 0.5 | [0; 1] | $0.7 \pm 0.1$ |
| $\sigma$ | 1 | [0,2] | 0.95 |

Table 3: Case study 2 (b): estimated parameter values for the Goutsias model of trascriptional regulation (Fig. 7).
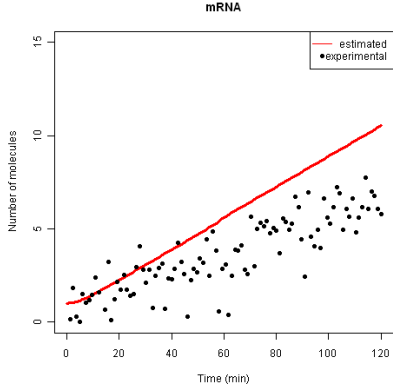
**mRNA**

Figure 6: simulations of case study 2. The actual and the estimated time behavior of the number of molecules of mRNA in the model network of Fig. 5.
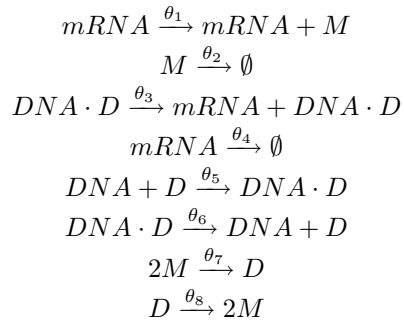
We also considered a more complex network model: the Goutsias model of gene transcription regulation [6, 13]. Figure 7 illustrates this model. The mRNA is translated into a protein monomer M that can dimerise. The dimer D, in turn, can bind to its DNA and acts as a transcription factor to auto-regulate its own mRNA production. Both mRNA and protein are degraded at constant rates. The set of reactions of this network is the following is also reported in Fig. 7. As in [13], we used this set of reactions to generate, with the Dizzy simulator (http://magnet.systemsbiology.net/software/Dizzy/), a synthetic dataset of the time series of the number of molecules for each component in the system. The dataset contains of 100 data point at the time resolution of 1.2 min. As initial values we used M = 2, D = 4, DNA = 2, and mRNA = 0, DNA·D = 0. All the reaction constants are in units of per seconds.



$$mRNA \xrightarrow{\theta_1} mRNA + M$$
$$M \xrightarrow{\theta_2} \emptyset$$
$$DNA \cdot D \xrightarrow{\theta_3} mRNA + DNA \cdot D$$
$$mRNA \xrightarrow{\theta_4} \emptyset$$
$$DNA + D \xrightarrow{\theta_5} DNA \cdot D$$
$$DNA \cdot D \xrightarrow{\theta_6} DNA + D$$
$$2M \xrightarrow{\theta_7} D$$
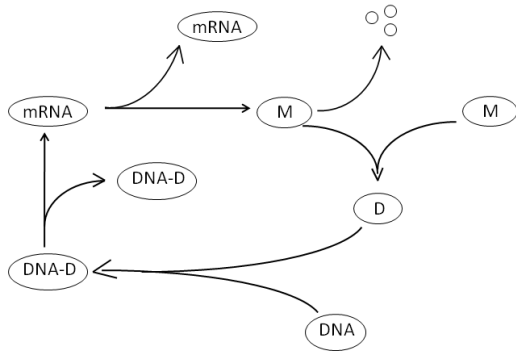$$D \xrightarrow{\theta_8} 2M$$

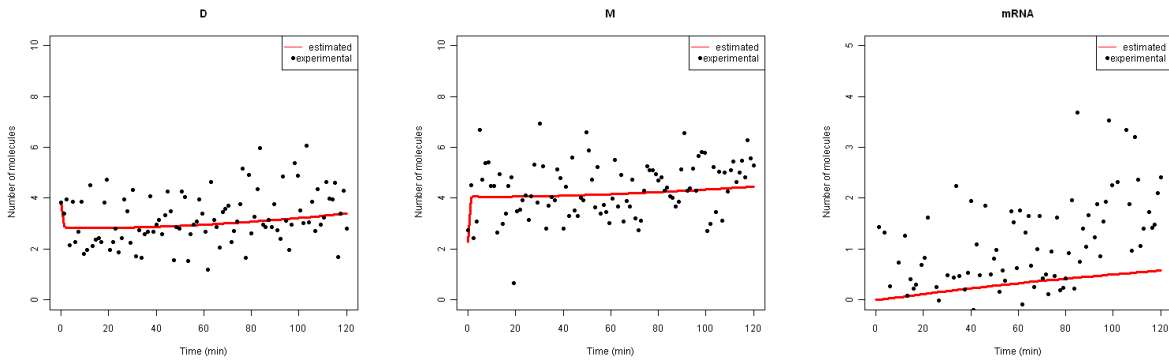Figure 7: Goutsias's model of gene transcription regulation.



Figure 8: estimated and experimental time behavior of the dimer (D), monomer (M), and mRNA.

## 4.3 Case study 3: the Lotka reactions

The set of coupled, autocatalytic reactions of Lotka are given in Fig. 9. We tested our procedure on this system, which, even it seems structurally quite simple, possesses remarkable dynamical properties. Although

the predator-prey interpretation of the Lotka reactions is a bit crude, it is helpful to visualize the dynamics of this system. The first reaction describes how a certain predator species $Y_2$ reproduces by feeding on a certain prey species $Y_1$ ; the second reaction describes how $Y_1$ reproduces by feeding on a certain foodstuff, which is assumed to be only insignificantly depleted thereby; and the third reaction describes the eventual demise of $Y_2$ through natural causes. The correct estimation of the rate constants in this model is extremely important, because the its dynamics is particularly sensitive to their changes: even very small differences in this values can determine the presence or the absence of the oscillatory behavior of the predators and preys. The correct estimation of the rate constants in this model is extremely important, because the its dynamics is particularly sensitive to their changes: even very small differences in this values can determine the presence or the absence of the oscillatory behavior of the predators and preys. We generated a synthetic dataset of time-course of the amounts of $X, Y_1, Y_2$ to use as experimental input data to KInfer with the following values of the rate coefficients $\theta_1 = 0.01$, $\theta_2 = 0.0001$, $\theta_3 = 10$, and the following initial amounts $Y_1 = Y_2 = 10^3$; $X = 10^5$, and $Z = 0$. as in [20]. We generated 100 data points at time steps of about 0.3.
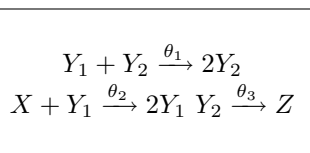
$$Y_1 + Y_2 \xrightarrow{\theta_1} 2Y_2$$
$$X + Y_1 \xrightarrow{\theta_2} 2Y_1 \; Y_2 \xrightarrow{\theta_3} Z$$

Figure 9: reactions of the Lotka model.

In this model the units of time and concentration/amount are not relevant for our purporses, since they are related to the specific organism/chemical we want to consider. The estimates obtained with Kinfer are reported in Table 4 and the comparison between estimated and experimental behavior is shoed in Fig. 10.

| Parameter | Actual value | Initial guesses | Estimated value |
|-----------|-------------|-----------------|-----------------|
| $\theta_1$ | 0.01 | [0; 0.002] | $0.0108 \pm 0.007$ |
| $\theta_2$ | 0.0001 | $[0; 10^{-4}]$ | $1.046 \times 10^{-4} \pm 10^{-6}$ |
| $\theta_3$ | 10 | [0; 12] | $11 \pm 1$ |
| $sigma$ | 1 | [0; 1] | 0.84 |

Table 4: case study 3: estimated parameter values for the Lotka reactions model.

The estimated behavior of $Y_1$ is oscillatory as the experimental one, but it has a smaller oscillation's amplitude. The period of the oscillations are almost the same instead. This proves the high sensitivity of this model to the slight variation of $\theta_2$.

## 4.4   Case study 4: neuron dynamics

The electrical properties of a segment of neuron membrane can be modeled by an equivalent circuit of the form shown in Fig. 11. In the equivalent circuit, current flow across the membrane has two major components, one associated with charging the membrane capacitance $C$ and one associated with the movement of specific types of ions across the membrane, and a non-specific leak current $I_L$. The ionic current is further subdivided into two main distinct components, a sodium current $I_{Na}$, a potassium current $I_K$. The fundamental equations describing the generation of action potentials and their spiking property had been established by Hodgkin and Huxley in early 1952 [7]. The Hodgking-Huxley equations (HH equations) comprise four highly non-linear differential equations, describing the dependence of neuron membrane potential $V$ on the flux of sodium and potassium ions. The HH equations model action potential or *spike generation* in the giant axon of the squid. Neocortical neurons in humans and other mammals are much more complex, however, as there are a total of at least 12 ion currents present in their membranes. To describe all of these currents in detail requires 16 coupled nonlinear differential equations. However, a few of these currents appear to contribute the majority of dynamical properties underlying neocortical firing patterns, so that these equations can be simplified to an excellent degree of approximation [18] to the form that is showed in Fig. 11. This model is intended to achieve sufficient simplification of the numerical calculations to permit its use in modeling small cortical neural networks. In the equation of this table, $C$ is the neuron membrane capacitance, $W$ is
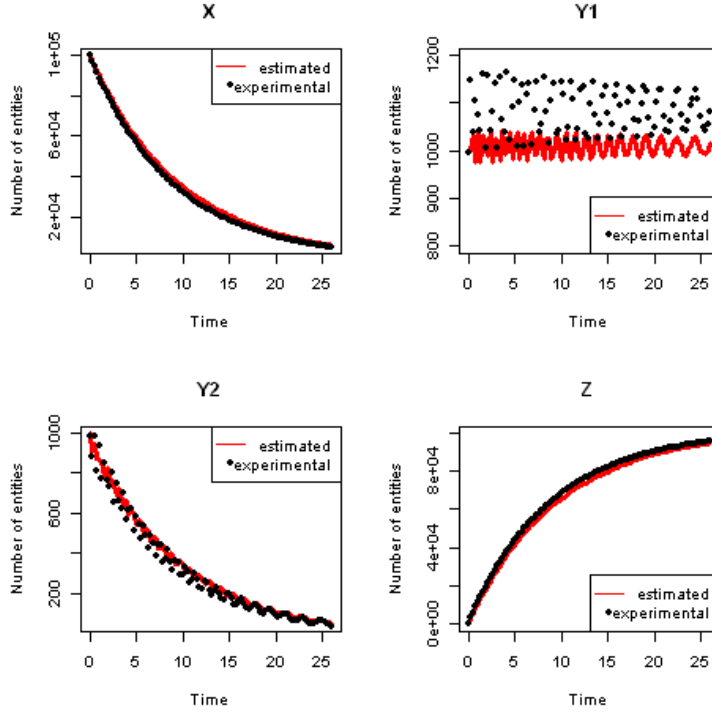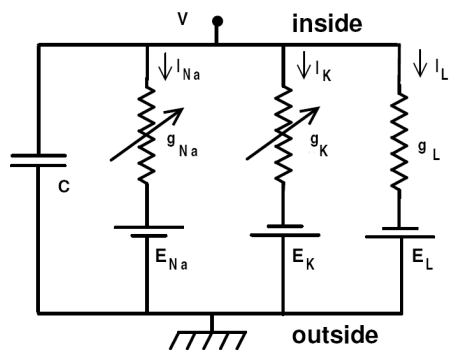
Figure 10: simulations of case study 3. Experimental and estimated behavior of Lotka model reagents $X, Y_1, Y_2, Z$.

the K$^+$ conductance mediating the recovery phase of the action membrane potential $V$, $g_{Na}$ is the electric conductance of Na$^+$, and $E_{Na}$ and $E_K$ are the equilibrium potentials of the Na+ and K$^+$ ions, respectively. Depending on values assigned to the parameters $I_{input}$, $\tau_W$ and $C$ a broad spectrum of neocortical activity patterns, including burst firing, can be simulated. Here only one will be emphasized: the spike frequency adaptation of neocortical regular spiking neurons in response to a constant stimulus $I_{input}$ and for real physiological values of $C$, $E_{Na}$ $E_K$, and $\tau_W$ (see Fig. 11).

Developing the equations in Fig. 11, we find that the non-linearities inherent in HH can be represented by cubic polynomials,as follows:

$$
\frac{dV}{dt} = -aV^3 - bV^2 + cV + dVW + eW + f, \quad \frac{dW}{dt} = gV - hW - i
$$
$$
a = -32/C, \ b = (-47.71 + 32.63)/C, \ c = (-18.81 + 47.71E_{Na})/C
$$
$$
d = 16/C, \ e = 16E_K)/C, \ f = (18.81E_{Na})/C + I_{in}/C
$$
$$
g = 1.35/\tau, \ h = 1/\tau, \ i = 1.03/\tau
$$

The Figs. 12 and 13 proves that the estimated parameters reproduces the experimental dynamics within the experimental uncertainties: experimental data are affected by a level of noise of $\sigma = 0.4$ that propagates - unchanged - to the estimated data. The discrepancies between actual and estimated value of the constants $c, d, e, f, i$ do not affect the good match between the actual and estimated dynamics. They are due to inaccuracies introduced by the discretization of the derivative $dV/dt$, that, anyway influence parameters for which the model is not sensitive.

$$\frac{dV}{dt} = \frac{1}{C}[-g_{Na}(V)(V - E_{Na}) - W(V - E_K) + I_{input}]$$

$$\frac{dW}{dt} = \frac{1}{\tau_W}(-W + G(V))$$

$$g_{Na}(V) = (18.81 + 47.71\ V)\ \text{nS}$$

$$G(V) = (1.03 + 1.35\ V)nS$$

$$C = 20\mu\text{F}/\text{cm}^2,\ E_{Na} = 55\ \text{mV},\ E_K = -92\ \text{mV}$$

$$\tau_W = 5\ \text{msec},\ I_{in} = 1.5\ \text{pA}$$

Figure 11: nn the left, the equivalent HH circuit for an electrically active membrane. The capacitance is due to the physiological bilayer separating the ions on the inside and the outside of the cell. The conductance of the $Na^+$ and $K^+$ currents are voltage dependent, as indicated by variable resistances. On the right, the reduced form Hodgkin-Huxley equations and the model parameters.

Using the physiologically reasonable values for $E$, $g$, $\tau$ and $g$, we reproduced typical experimental data of the action potential [18, 7]. In order to obtain real experimental data from the HH model we rescaled $V$ in the following way: $V \leftarrow (V + |\min(V)|) \cdot 24)$ and we used these rescaled data as input to KInfer. Table 5 reports the estimated values for the parameters $a, b, \ldots, l$ that are combinations of the $E$, $g$, $\tau$ and $g$, as in Fig. 11 (their values are accordingly rescaled to the rescaling of $V$). In this table we did not report the bound for the initial guesses, as they are extremely close to the actual expected values.

| Parameter | Actual value | Estimated value |
|:---:|:---:|:---:|
| $a$ | 1.6315 | 1.6225 |
| $b$ | 1.4881 | 1.4138 |
| $c$ | 0.3715 | 0.1414 |
| $d$ | 0.8000 | 0.4580 |
| $e$ | 0.7360 | 0.4319 |
| $f$ | 0.5923 | 0.3247 |
| $g$ | 0.2700 | 0.3000 |
| $h$ | 0.2000 | 0.2246 |
| $i$ | 0.2060 | 0.05 |
| $\sigma$ | 0.4 | 0.4 |

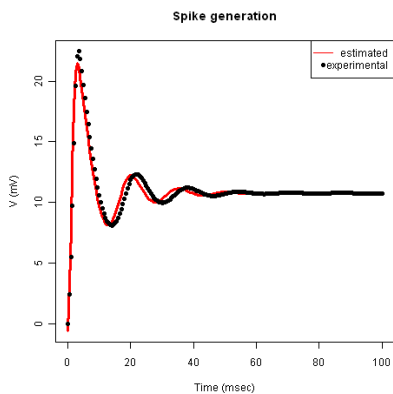Table 5: Case study 4: estimated parameter values for the HH model.



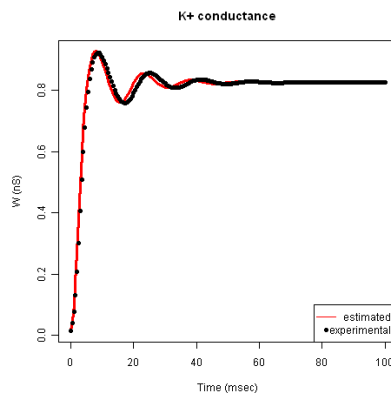Figure 12: estimated versus experimental behavior of spike generation.



Figure 13: estimated versus experimental behavior of potassium conductance.

14

# 5    Conclusions

In this article, we presented a novel method for the estimation of reaction parameters and noise strength from time series of molecules counts or concentrations observed with error. We have shown that our procedure converges to the expected solutions within the bounds of the experimental errors that propagates from concentration measurements to the kinetic rate constants. The results confirm that the validity of the procedure and of the discretized model of mass action law for the rate equation. Moreover, some important features missing from the existing methods for parameter inference are present in our method. The first is the implementation of a procedure, which automates the computation of the initial guesses of the parameters. In this way, the user is not forced to insert any *a priori* knowledge about the system, that often is quite hard to find, and, at the same time, the method is equipped with a rigorous procedure referring only to the experimental concentration measurements to identify a region of the parameter space where the optimization of the probability density function takes place. The second feature is the implementation of the experimental error propagation. The evaluation of the experimental uncertainty on the rate constants estimates is particularly useful if the procedure of parameter inference is incorporated in projects of experimental design. The size of the errors on the kinetic constants is indicative of the optimality of the experimental setup. Thus, any procedure devoted to the reduction of this error is definitely part of a methodology aiming to optimize the design of the experimental configuration.

# References

[1] Boys, R. J., Wilkinson, D. J., & Kirkwood, T. B. (2008). Bayesian inference for a discretely observed stochastic kinetic model. Statistics and Computing, Springer Netherlands.

[2] Chou, I.-C., Martens, H., & Voit, E. O. (2006). Parameter estimation in biochemical systems models with alternating regression. Theoretical Biology and Medical Modelling , 3, 25.

[3] Goldberg, D. E. (1989). Genetic algorithms in search, optimization and machine learning. Addison-Wesley, Massachusetts.

[4] Golding, I., Paulsson, & Zawilski, S. M. (2005). Real-time kinetics of gene activity in individual bacteria. Cell , 123, 1025-1036.

[5] Golightly, A., & Wilkinson, D. J. (2008). Bayesian inference for nonlinear multivariate diffusion models observed with error. Computational statistics and data analysis , 52 (3), 1674-1693.

[6] Goutsias, J. (2006). A hidden Markov model for transcriptional regulation in single cells, IEEE/ACM Trans. Comput. Biol. Bioinform. , 3, 57-71.

[7] Hodgkin A. L. & Huxley A. F. (1952) A quantitavive description of membrane current and its application to conduction and exitation in nerve, The Journal of Physiology.

[8] Hoops, S., Sahle, S., Gauges, R., Lee, C., Pahle, J., Simus, N., et al. (2006). COPASI - a COmplex PAthway SImulator. Bioinformatics , 22, 3067-3074.

[9] Self citation (2007). A new method for inferring rate coefficients from experimental time-consecutive measurement of reactant concentrations. Int. Conference on Systems Biology. Long Beach.

[10] Self citation (2008). Calibration of biochemical network models. Technical Report The Microsoft Research - University of Trento Centre for Computational and Systems Biology n. 16/2008.

[11] Moles, G. C., Mendes, P., & Banga, J. R. (2003). Parameter estimation in biochemical pathways: a comparison of global optimization methods. Genome Res. (13), 2467-2474.

[12] Polisetty, P. K., Voit, E. O., & Gatzke, E. P. (2006). Identification of metabolic system paramters using global optimization methods. Theoreteical Biology and Medical Modelling , 3 (4).

[13] Reinker, S., Altman, R. M., & Timmer, J. (2006). Parameter estimation in stochastic biochemical reactions. 153 (4).

[14] Rodrigez-Fernandez, M., Mendes, P., & Banga, J. (2006). A hybrid approach for efficient and robust parameter estimation in biochemical pathways. BioSystems , 83: 248-265.

[15] Sugimoto, M., Kikuchi, S., & Tomita, M. (2005). Reverse engineering of biochemical equations from time-course data by means of genetic programming. BioSystems , 80: 155-164.

[16] Tian, T., Xu, S., & Burrage, K. (2007). Simulated maximum likelihood method for estimating kinetic rates in gene expression. Bioinformatics , 23(1), 84-91.

[17] Wilkinson, D. J. (2007). Bayesian methods in bioinformatics and computational systems biology. Briefings in bioinformatics , 1-8.

[18] Wilson H. R. (1999). Nonlinear dynamics of neurons and network in vision, Procs. of the $39^{th}$ Conference on Decision & Control, Phoenix, USA.

[19] Zwolak, J., PET - Parameter Estimation Toolkit. Retrieved 2007, Software's home page: http://mpf.biol.vt.edu/pet/contact.php

[20] Gillespie, D. T. (1977). Exact stochastic simulation of coupled chemical reactions. J. of Physical Chemistry, 81 (25).

[21] BioBayes: Bayesian Inference for Systems Biology ; http://www.dcs.gla.ac.uk/BioBayes/.