

Data Assimilation for Systems and Mathematical Biology

Benjamin Engelhardt^{1,2}, Dominik Kahl³, Andreas Weber⁴, Maik Kschischo^{3,*}

March 27, 2018

Abstract

Mathematical models are increasingly used as a tool to deal with the tremendous complexity of biological systems. Data Assimilation, defined as the process of combining models with experimental observations, is a key step in order to better align the model outputs with reality. Despite new and improved experimental techniques, it is usually impossible to directly observe all the states of a biological system. This renders Data Assimilation an inverse problem requiring sophisticated mathematical and statistical techniques and the systematic integration of prior knowledge or assumptions.

In this article, we will highlight how Data Assimilation is key to dealing with incomplete and uncertain biological information. We will survey the basic concepts and methodological approaches relevant to systems biology and indicate similarities and differences with Data Assimilation problems faced in the environmental and geosciences. In addition, we will illustrate the Dynamic Elastic-Net, a recent data driven approach to dealing with incomplete models and structural model errors.

1 Introduction

A central goal of the life sciences is to understand, predict or manipulate the dynamics of biological systems. For example, ecologists record the number of individuals of a group or species in a certain area to reveal the inner workings of dynamical ecosystems. Molecular and cell biologists measure the abundance of proteins and other gene products as a function of time in order to better characterise molecular pathways. The central theme of mathematical and computational systems biology is to combine these data with mathematical models and to use computer simulations to better understand the system, to predict its behaviour after certain interventions and to design controls to manipulate the system in a desired way.

One of the greatest challenges of these combined modelling and experimentation endeavours is Data Assimilation (DA), defined here as the process of combining the mathematical model with the observed data. An important DA problem is the estimation of dynamic state variable in the rather typical situation that not all the variables of a state space model can directly be experimentally accessed. This state observer problem is well known in control engineering and a great part of research has been devoted to deriving conditions under which the state can be reconstructed from the outputs (observability) and to design state observers, i.e. artificial dynamical systems synchronising their state to the dynamics of the true system using only the available output data [TOCITE]. Measurement noise makes the state estimation problem even harder. As in other areas of science and engineering, measurement noise originates from shortcomings of the measurement

¹ Rheinische Friedrich-Wilhelms-Universität Bonn, Algorithmic Bioinformatics, Bonn, Germany

² Current Address: AbbVie Deutschland GmbH & Co. KG, Ludwigshafen, Germany

³ University of Applied Sciences Koblenz, RheinAhrCampus, Department of Mathematics and Technology, Remagen, Germany

⁴ Rheinische Friedrich-Wilhelms-Universität Bonn, Institute for Computer Science, Bonn, Germany

*Corresponding author: kschischo@rheinahrcampus.de

device or the experimental procedures. If the system states itself can be considered to be deterministic, the state estimation problem reduces to the estimation of parameters and initial conditions. The challenges of parameter estimation have extensively been studied in the context of systems biology [TOCITE].

In Fig.1 we illustrate the state estimation problem using a deterministic dynamical model for the regulation of the Epo receptor [1]. At the moment, we will refrain from discussing the model in detail and just focus on the principled state estimation problem. The model has six state variables $\mathbf{x} = (x_1, \dots, x_6)^T$ and their dynamics $\mathbf{x}(t)$ is described by a system of coupled ordinary differential equations (ODEs). These interactions can be displayed graphically, see section 2.1 for details. The state variables represent the concentrations of different biomolecules, which can not directly be measured. Instead, the measurable outputs $\mathbf{y} = (y_1, y_2, y_3)^T$ of the system are scaled linear combinations of the states. The DA task involves the estimation of the parameters and the initial state from noisy time course observations $\mathbf{y}(t_k)$ at certain measurement times t_k . Given that the system is deterministic and has a unique solution, the parameters and the initial condition provide full information about the dynamic state $\mathbf{x}(t)$.

A fundamentally different source of noise in living systems is related to biological diversity. No two apparently alike cells in an organ or no two individuals of the same species are in fact completely identical or react in exactly the same way to an environmental stimulus. This variation is the basis for natural selection, but it also poses additional challenges for modelling and DA. In state space models, the diversity can be represented by stochastic dynamical systems, where the state variables are targeted by noise. This is sometimes called systems noise. Other sources of systems noise can be environmental fluctuations not encompassed by the model. Thus, both intrinsic and extrinsic noise processes lead to stochastic dynamics and the task of state estimation is then expressed as the problem of inferring the probability distribution of the states given the observed data.

The majority of DA algorithms is derived from the assumption that the governing equations of the system are known. However, in biology we have only partial information about all the interactions and processes in a real living system and the decision of which terms in a model can be neglected is often not easy to make. In addition, biological systems are open and exchange matter, energy and information with their environment. Thus, the perfect model assumption is usually unrealistic and structural model errors can not be represented by zero mean noise processes. DA algorithms dealing with structural model error are still in its infancy, but progress to deal with model errors and incomplete information is urgently required.

In our view, DA is a key technique in all areas of quantitative biology, although the most recent methodological progress has been largely driven by applications in other fields, in particular the atmospheric and oceanographic sciences. Our motivation is to review the most important concepts of DA in the context of quantitative biology and to introduce researchers in the life sciences to recent developments and open research questions relevant to their modelling work.

2 The Data Assimilation problem

2.1 Continuous time state space models

We assume that the state of a biological system can be characterised by a vector of state variables $\mathbf{x}(t) = (x_1(t), \dots, x_n(t))^T$ at time t . One example is the state of a biochemical reaction network, which is given by the concentrations x_i of all the reacting substances i . Other examples are the state of a cellular population, which can be characterised by the abundance x_i of each cell type i or the number of individuals x_i of a species i in a population dynamics model. For continuous time $t \in \mathbb{R}_{\geq 0}$, the dynamical model is formulated as a system of

ordinary differential equations

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)) \quad (1a)$$

$$\mathbf{y}(t) = \mathbf{h}(\mathbf{x}(t)). \quad (1b)$$

$$\mathbf{x}(t_0) = \mathbf{x}_0, \quad (1c)$$

where the function $\mathbf{f} : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}^n$ with $\mathcal{X} \subset \mathbb{R}^n$, $\mathcal{U} \subset \mathbb{R}^q$ encodes the effect of the interacting state variables and the known input $\mathbf{u}(t) \in \mathcal{U}$ on the time derivative $\dot{\mathbf{x}}(t)$. Note, that one can also arrive at such an ordinary differential equation model (1) from spatio-temporal models by discretising in space. Typically, it is not possible to obtain direct measurements for all the state variables. Instead, one observes the output variables $\mathbf{y}(t) = (y_1(t), \dots, y_m(t))^T$, which are related to the states in (1) via the measurement function \mathbf{h} . Usually, we will set $t_0 = 0$ in the initial condition (1c).

To illustrate different aspects of DA we will use a model for the information processing at the erythropoietin (Epo) receptor (EpoR) as a running example [1]. This model is illustrated in Fig. 1 together with its influence graph, which has a node for each state variable x_i and an directed edge from x_i to x_j whenever $\frac{\partial f_j}{\partial x_i} \neq 0$. In addition, the effect of the state to the outputs is visualised by drawing a dashed arrow from x_i to y_l , if $\frac{\partial g_l}{\partial x_i} \neq 0$.

2.2 The Filtering, smoothing and prediction tasks

In reality, one observes the output \mathbf{y} of a state space model at discrete measurement time points $t_1 < t_2 < \dots < t_T$. The observations $\mathbf{y}^o(t_k)$ are in addition corrupted by measurement noise $\mathbf{v}(t_k)$ and thus

$$\mathbf{y}^o(t_k) = \mathbf{h}(\mathbf{x}(t_k)) + \mathbf{v}(t_k) \quad (2)$$

describes the sequence of measurements, cf. (1b). We assume that the measurements devices are unbiased and thus the noise $\mathbf{v}(t_k)$ is a stochastic process with expectation zero: $E(\mathbf{v}(t_k)) = \mathbf{0}_{\mathbb{R}^m}$ for all times $k = 1, 2, \dots, n$.

Assume that we have observed the output of our dynamical system and collected an $m \times T$ data matrix $\mathbf{y}_{1:T}^o = (\mathbf{y}^o(t_1), \mathbf{y}^o(t_2), \dots, \mathbf{y}^o(t_T)) \in \mathbb{R}^{m \times T}$. The goal

of state estimation is to combine the data with the model (1) in order to infer the state $\mathbf{x}(t)$ at time t . Depending on t , the task of estimating the state $\mathbf{x}(t)$ given the data $\mathbf{y}_{1:T}^o$ is called

(a) *smoothing* if $t \in [t_1, t_T]$

(b) *filtering* if $t = t_T$

(c) *prediction* if $t > t_T$

Thus, smoothing is an offline task, where we use all the observations to infer the state, whereas filtering refers to the online task of estimating the state immediately after observing $\mathbf{y}(t_T)$. Predictions are estimates of the future behaviour beyond the data assimilation window $[t_1, t_T]$.

3 Discrete time state space model

The continuous-time dynamical system (1a) can also be considered at discrete time points via the map

$$\begin{aligned} \phi_s(\mathbf{x}(t_0)) &= \mathbf{x}(t_0 + s) \\ &= \mathbf{x}(t_0) + \int_{t_0}^{t_0+s} \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)) dt. \end{aligned} \quad (3)$$

The resulting discrete-time dynamic system

$$\mathbf{x}_{k+1} = \phi_t(\mathbf{x}_k), \quad (4)$$

describes the dynamics of $\mathbf{x}_k = \mathbf{x}(ks)$. Here, we have suppressed the known input \mathbf{u} in the notation. An advantage of the discrete-time description is that the technicalities of stochastic differential equations can be avoided, when we consider systems with state noise. Thus, in the following derivations we will consider the general model

$$\mathbf{x}_{k+1} = \phi_t(\mathbf{x}_k) + \mathbf{w}_k \quad (5a)$$

$$\mathbf{y}_k = \mathbf{h}(\mathbf{x}_k(t)) + \mathbf{v}_k \quad (5b)$$

where \mathbf{w}_k and \mathbf{v}_k describe the state and observation noise, respectively. These stochastic processes are assumed to be independent and

identically distributed (i.i.d) sequences of random variables densities p_w and p_v . We use the notation

$$\mathbf{w}_k \sim p_w \quad \text{and} \quad \mathbf{v}_k \sim p_v \quad \forall k \in \mathbb{N}_0. \quad (5c)$$

We assume zero expectation $E(w_k) = 0$ and $E(v_k) = 0$. In addition, the initial state \mathbf{x}_0 is often unknown and thus considered a random variable $\mathbf{x}_0 \sim p(\mathbf{x}_0)$.

The stochastic dynamics (5) can be reformulated as a probabilistic state space model. The probability density $p(\mathbf{x}_{0:T})$ of a state sequence $\mathbf{x}_{0:T} = (\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T)$ and the conditional density of an observation sequence $\mathbf{y}_{1:T} = (\mathbf{y}_1, \dots, \mathbf{y}_T)$ given the state sequence can be recursively computed as

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_0) \prod_{k=1}^T p(\mathbf{x}_k | \mathbf{x}_{k-1}) \quad (6)$$

$$p(\mathbf{y}_{1:T} | \mathbf{x}_{0:T}) = \prod_{k=1}^T p(\mathbf{y}_k | \mathbf{x}_k). \quad (7)$$

Here, (6) indicates that the state space model (5a) defines a Markov process of first order, where the state \mathbf{x}_k at time step k depends only on the previous state \mathbf{x}_{k-1} . The state transition probabilities are given as

$$p(\mathbf{x}_{k+1} | \mathbf{x}_k) = p_w(\mathbf{x}_{k+1} - \phi_t(\mathbf{x}_k)). \quad (8)$$

The product in (7) expresses the conditional independence of the outputs \mathbf{y}_k at different time points, given the state \mathbf{x}_k . By comparison with (5b) we have

$$p(\mathbf{y}_k | \mathbf{x}_k) = p_v(\mathbf{y}_k - \mathbf{h}(\mathbf{x}_k)). \quad (9)$$

In summary, we see that the discrete time state space model (5) together with specifications of the noise distributions (5c) is equivalent to specifying a state transition probability $p(\mathbf{x}_{k+1} | \mathbf{x}_k)$ and an output distribution $p(\mathbf{y}_k | \mathbf{x}_k)$. Readers familiar with sequence analysis (TOCITE-DURBIN) recognise the structure of a Hidden Markov Model (HMM) with hidden state \mathbf{x}_k and observations \mathbf{y}_k . However, here we consider continuous state spaces \mathcal{X} instead of the discrete state spaces used for protein or nucleotide sequence analysis.

4 The Smoothing and the filtering distributions

For DA, we want to infer the states given the observations. The posterior distribution of the states follows from Bayes' theorem as

$$p(\mathbf{x}_{0:t} | \mathbf{y}_{1:T}) = \frac{p(\mathbf{y}_{1:T} | \mathbf{x}_{0:t}) p(\mathbf{x}_{0:T})}{p(\mathbf{y}_{1:T})} \quad (10)$$

$$= \frac{p(\mathbf{x}_0) \prod_{k=1}^T p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{y}_k | \mathbf{x}_k)}{p(\mathbf{y}_{1:T})} \quad (11)$$

with the marginal output distribution

$$p(\mathbf{y}_{1:T}) = \int p(\mathbf{y}_{1:T} | \mathbf{x}_{0:t}) p(\mathbf{x}_{0:T}) d\mathbf{x}_0 \dots d\mathbf{x}_T. \quad (12)$$

If we insert the observed data $\mathbf{y}_{1:T} = \mathbf{y}_{1:T}^o$ into the posterior smoothing distribution, then the posterior smoothing distribution (10) provides maximum information about the states. Inferring this distribution is called the *smoothing problem*. As we will see, the integral in (12) is often difficult or impossible to calculate exactly. Smoothing algorithms approximate the smoothing distribution or certain of its statistical moments (10).

Filtering refers to the online task of estimating \mathbf{x}_k from data $\mathbf{y}_{1:k}^o$ accumulated up to time k . Thus, we are interested in the filtering distribution $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ for time steps $k = 1, 2, \dots, T$. This filtering distribution is updated from $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ to $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k+1})$ via two steps, *prediction* and *analysis*. Prediction refers to the mapping from $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ to $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})$, i.e. prediction of the next state given the output data up to time point k . The corresponding distribution $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})$ is often called the *background* distribution. Analysis considers the mapping from the background $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})$ to the *analysis* distribution $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k+1})$ to incorporate the next observation \mathbf{y}_{k+1} into the state estimate.

Prediction: $p(\mathbf{x}_k | \mathbf{y}_{1:k}) \rightarrow p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})$

The prediction step is based on the Chapman-Kolmogorov equation

$$p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k}) = \int_{\mathcal{X}} p(\mathbf{x}_{k+1} | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k}) d\mathbf{x}_k. \quad (13)$$

Here, we have used $p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{y}_k) = p(\mathbf{x}_{k+1}|\mathbf{x}_k)$, which means that the information about the perfect state \mathbf{x}_k cannot be improved by the noisy measurement \mathbf{y}_k . In the special case of a deterministic state space model, i.e. in the absence of process noise ($\mathbf{w}_k = 0$ for all $k \in \mathbb{N}_0$ in (5a)), the state transition probability density degenerates to a delta distribution $p(\mathbf{x}_{k+1}|\mathbf{x}_k) = \delta(\mathbf{x}_{k+1} - \phi_t(\mathbf{x}_k))$ and the prediction step corresponds to a single iteration of the deterministic discrete time state space model.

Analysis: $p(\mathbf{x}_{k+1}|\mathbf{y}_{1:k}) \rightarrow p(\mathbf{x}_{k+1}|\mathbf{y}_{1:k+1})$
The analysis step is based on an application of Bayes' theorem

$$\begin{aligned} p(\mathbf{x}_{k+1}|\mathbf{y}_{1:k+1}) &= p(\mathbf{x}_{k+1}|\mathbf{y}_{k+1}, \mathbf{y}_{1:k}) \\ &= \frac{p(\mathbf{y}_{k+1}|\mathbf{x}_{k+1}, \mathbf{y}_{1:k}) p(\mathbf{x}_{k+1}|\mathbf{y}_{1:k})}{p(\mathbf{y}_{k+1}|\mathbf{y}_{1:k})} \\ &= \frac{p(\mathbf{y}_{k+1}|\mathbf{x}_{k+1}) p(\mathbf{x}_{k+1}|\mathbf{y}_{1:k})}{p(\mathbf{y}_{k+1}|\mathbf{y}_{1:k})}, \end{aligned} \quad (14)$$

where in the last line we have used once again the conditional independence of the observation from previous observations given the current state. Together, the two equations (13, 14) are approximated by different filtering algorithms to recursively estimate the filtering distribution $p(\mathbf{x}_k|\mathbf{y}_{1:k})$. This recursive approach is often called *sequential data assimilation*.

5 Variational smoothing and filtering

Inferring the smoothing and filtering posterior distributions can be challenging and computationally expensive. In some cases, it is sufficient to summarise these distributions by point estimates. Variational methods for smoothing and filtering seek a maximum a posteriori estimator (MAP estimator), i.e. they locate peaks of the smoothing or filtering distribution, respectively, by solving an optimisation problem.

The posterior smoothing distribution (10) can

be rewritten as

$$\begin{aligned} p(\mathbf{x}_{0:t}|\mathbf{y}_{1:T}) &= e^{-\Psi(\mathbf{x}_{0:t}|\mathbf{y}_{1:T})} \\ &\propto e^{-E_0(\mathbf{x}_0) + \sum_{k=1}^T [E(\mathbf{x}_k|\mathbf{x}_{k-1}) + U(\mathbf{y}_k|\mathbf{x}_k)]}, \end{aligned} \quad (15)$$

where we have assumed that the state transition and the output probabilities are nonzero everywhere on their domain and

$$\begin{aligned} p(\mathbf{x}_0) &\propto \exp(-E_0(\mathbf{x}_0)) \\ p(\mathbf{x}_k|\mathbf{x}_{k-1}) &\propto \exp(-E(\mathbf{x}_k|\mathbf{x}_{k-1})) \\ p(\mathbf{y}_k|\mathbf{x}_k) &\propto \exp(-U(\mathbf{y}_k|\mathbf{x}_k)). \end{aligned} \quad (16)$$

The constant term originating from the denominator in (10) does not depend on the state sequence and is omitted in the proportionality (15). The MAP estimator for the smoothing problem is obtained from minimising $\Psi(\mathbf{x}_{0:t}|\mathbf{y}_{1:T}^o)$ over all possible state sequences $\mathbf{x}_{0:t}$ for given data $\mathbf{y}_{1:T}^o$. The constants of proportionality in (16) are the normalisation constants independent of the arguments of the three "energies" E_0 , E and U .

We illustrate this minimisation for the important case that the process noise \mathbf{w}_k , the measurement noise \mathbf{v}_k and the initial state \mathbf{x}_0 are iid Gaussian random variables, i.e. $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}_{\mathbb{R}^n}, Q)$, $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}_{\mathbb{R}^m}, R)$ and $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}_{\mathbb{R}^n}, S)$ with zero mean and covariance matrix Q , R and S respectively. Then*,

$$E_0(\mathbf{x}_0) = \frac{1}{2} \|\mathbf{x}_0\|_S^2 = \quad (17)$$

$$E(\mathbf{x}_k|\mathbf{x}_{k-1}) = \frac{1}{2} \|\mathbf{x}_{k+1} - \phi_t(\mathbf{x}_k)\|_Q^2 \quad (18)$$

$$U(\mathbf{y}_k|\mathbf{x}_k) = \frac{1}{2} \|\mathbf{y}_k - \mathbf{h}(\mathbf{x}_k)\|_R^2 \quad (19)$$

6 Linear models: The Kalman filter and smoother

BIS HIER

*Here we have used the notation $\mathbf{a}^T B \mathbf{a} = \|\mathbf{a}\|_B^2$ for a column vector $\mathbf{a} \in \mathbb{R}^n$ and a square matrix $B \in \mathbb{R}^{n \times n}$.

The variances and covariances between the different output measurements are collected in the covariance matrix $R_s = \text{Cov}(\mathbf{v}_s) = \text{cov}(v_k(t_s), v_l(t_s))_{k,l=1,\dots,m}$ at time t_s . In many cases the measurement noise can be assumed to be Gaussian $\mathbf{y} \sim \mathcal{N}(\mathbf{0}_{\mathbb{R}^m}, R)$ and is thus completely characterised by R . Non-gaussian measurements can often be transformed to a Gaussian distribution. In this review we assume Gaussian measurement noise, although other distributions of the measurement noise can also be considered, but often at the price of more complicated DA equations.

- State observation
- discrete time and noise

7 Monte Carlo algorithms for DA

8 Structural model error

9 Outlook

A Appendix name

References

- [1] Becker V, Schilling M, Bachmann J, Baumann U, Raue A, Maiwald T, et al. Covering a Broad Dynamic Range: Information Processing at the Erythropoietin Receptor. *Science*. 2010 Jun;328(5984):1404–1408. Available from: <http://www.sciencemag.org/cgi/doi/10.1126/science.1184913>.

bla

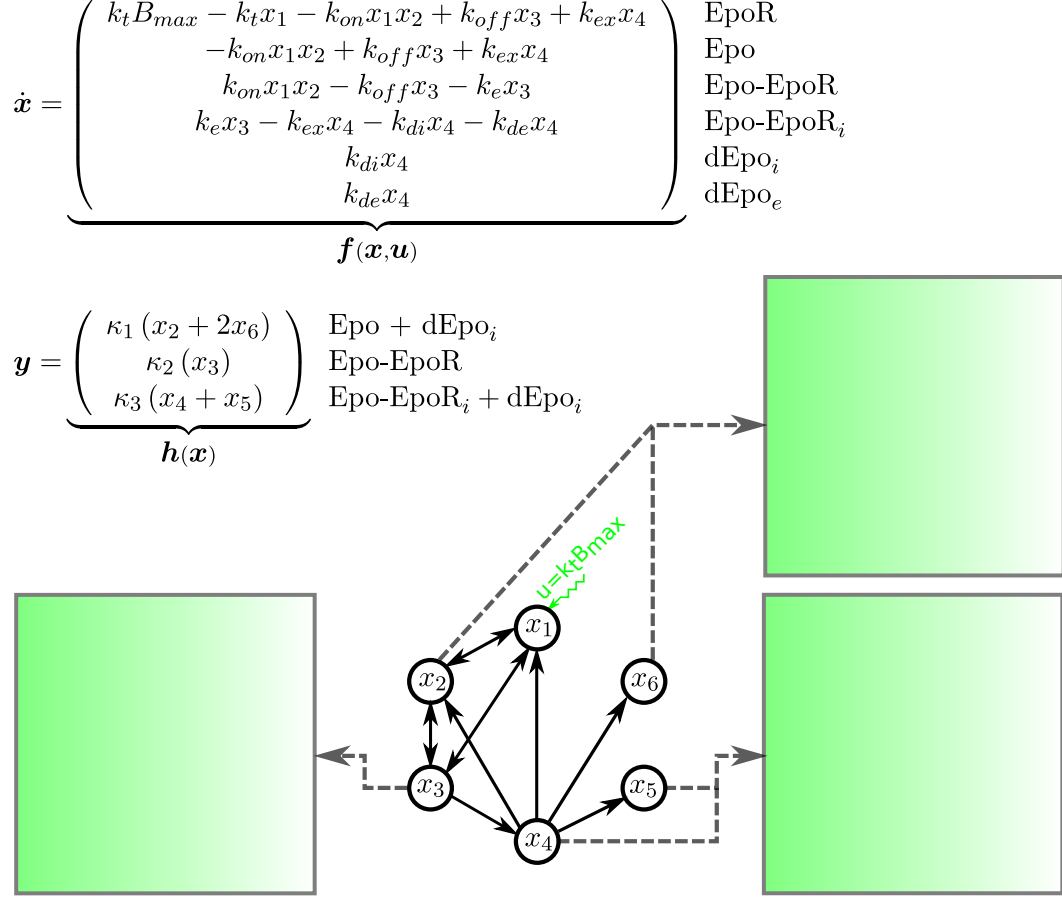


Figure 1: The Epo receptor regulation model [1] used as a running example. The state $\mathbf{x} = (x_1, \dots, x_6)^T$ of this model is given by the concentrations of the Epo receptor (x_1) on the cell surface which can bind to Epo (x_2) and build the ligand-receptor complex (x_3). This complex is able to activate subsequent signaling cascades, e.g. the JAK-STAT signaling pathway. In addition the ligand-receptor complex can be internalized (x_4) and dissociate from Epo which is then degraded (x_5) and transported to the extracellular space (x_6). The complex regulation of this receptor is characterized by receptor mobilization, turnover and recycling. The rate constants correspond to (i) receptor turnover (k_t), (ii) ligand-receptor binding (k_{on}) or dissociation (k_{off}), (iii) ligand-induced endocytosis (k_e), (iv) recycling (k_{ex}) and (v) internal (k_{di}) or external (k_{de}) degradation of Epo. Only the the Epo concentration in medium (y_1), on surface (y_2) and in cells (y_3) can be measured up to some scaling parameters κ_j , $j \in \{1, 2, 3\}$. The influence graph illustrates the interactions between the state variables (black arrows) and the measured outputs y_k , $k \in \{1, 2, 3\}$ (dashed arrows).