

Fast derivatives of likelihood functionals for ODE based models using
adjoint-state method

Peer-reviewed author version

Melicher, Valdemar; HABER, Tom & Vanroose, Wim (2017) Fast derivatives of
likelihood functionals for ODE based models using adjoint-state method. In:
COMPUTATIONAL STATISTICS, 32(4), p. 1621-1643.

DOI: 10.1007/s00180-017-0765-8

Handle: <http://hdl.handle.net/1942/25461>

Fast derivatives of likelihood functionals for ODE based models using adjoint-state method

Valdemar Melicher¹ · Tom Haber² · Wim Vanroose¹

Received: date / Accepted: date

Abstract We consider time series data modeled by ordinary differential equations (ODEs), widespread models in physics, chemistry, biology and science in general. The sensitivity analysis of such dynamical systems usually requires calculation of various derivatives with respect to the model parameters.

We employ the *adjoint state method* (ASM) for efficient computation of the first and the second derivatives of likelihood functionals constrained by ODEs with respect to the parameters of the underlying ODE model. Essentially, the gradient can be computed with a cost (measured by model evaluations) that is independent of the number of the ODE model parameters and the Hessian with a linear cost in the number of the parameters instead of the quadratic one. The sensitivity analysis becomes feasible even if the parametric space is high-dimensional.

The main contributions are derivation and rigorous analysis of the ASM in the statistical context, when the discrete data are coupled with the continuous ODE model. Further, we present a highly optimized implementation of the results and its benchmarks on a number of problems.

The results are directly applicable in (e.g.) maximum-likelihood estimation or Bayesian sampling of ODE based statistical models, allowing for faster, more stable estimation of parameters of the underlying ODE model.

Keywords: Sensitivity Analysis, Ordinary Differential Equations, Gradient, Hessian, Statistical Computing, Mathematical Statistics, Algorithm

The work of the first two authors was supported by IWT O&O project 130406 – ExaScience Life HPC.

Valdemar Melicher
E-mail: Valdemar.Melicher@UAntwerpen.be

¹ Department of Mathematics and Computer Science, University of Antwerp, Middelheimlaan 1, 2020 Antwerp, Belgium

² Expertise Centre for Digital Media, Hasselt University, Wetenschapspark 2, 3590 Diepenbeek, Belgium

1 Introduction

We consider time series vector data $\mathbf{y}_i \in \mathbb{R}^n$ for $i = 1, \dots, N$, where n is the dimension of the observation space and N is the number of corresponding measurements times t_i in the interval $I := [0, T]$ with some positive final time $T > 0$. In many scientific fields, the underlying structural model for such data is very often an initial-value problem of the following type:

$$\begin{aligned} d_t \mathbf{u} &= \mathbf{f}(t, \mathbf{u}, \boldsymbol{\phi}), \quad t \in [0, T], \\ \mathbf{u}(0) &= \mathbf{u}_0(\boldsymbol{\phi}), \end{aligned} \quad (1)$$

where \mathbf{u}_0 is the initial condition, dependent only on the parameter vector $\boldsymbol{\phi} \in \mathbb{R}^p$. In general non-linear r.h.s. \mathbf{f} of the governing equation represents the time derivative of the model variable $\mathbf{u}(t)$. It depends on the current time t , the model parameters $\boldsymbol{\phi}$ and the current values of $\mathbf{u} \in \mathbb{R}^m$.

The predictor $\hat{\mathbf{y}}$ of the data \mathbf{y} is a result of integration of the dynamical system (1) and a possible subsequent post-processing, for example aggregation. This can be expressed in mathematical terms as $\hat{\mathbf{y}} = \mathcal{P}(\mathbf{u}(t, \boldsymbol{\phi})) =: \mathbf{g}(t, \boldsymbol{\phi})$, where $\mathcal{P} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is the post-processing operator relating the solution \mathbf{u} to data.

The main aim of this paper is to efficiently compute the first and the second derivatives of negative log-likelihood functionals of the following form

$$l(\boldsymbol{\phi}) = \sum_i d(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi})), \quad (2)$$

with respect to $\boldsymbol{\phi}$. Here $d : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, \infty)$ is a sufficiently smooth distance function (metric) on \mathbb{R}^n . Equation 2 measures the fidelity between the model and the data.

1.1 Motivation

The most prominent example of negative log-likelihood functional (2) is obtained for error model

$$\mathbf{y}_i = \mathbf{g}(t_i, \boldsymbol{\phi}) + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim_{i.i.d.} \mathcal{N}(0, \boldsymbol{\Sigma}), \quad (3)$$

i.e. the residual errors $\boldsymbol{\varepsilon}_i$ are independent and identically distributed normal random variables with zero mean and residual covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$. Then

$$d(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi})) := \frac{1}{2} (\mathbf{y}_i - \mathbf{g}(t_i, \boldsymbol{\phi}))^t \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \mathbf{g}(t_i, \boldsymbol{\phi})) \quad (4)$$

and we are interested in the derivatives of

$$l(\boldsymbol{\phi}) := -\log p(\mathbf{y}|\boldsymbol{\phi}) \propto \sum_i d(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi})). \quad (5)$$

The gradient or the Hessian of such a negative log-likelihood are required quite often in statistics in various contexts. Let us supply a few examples. First, Laplace's method (approximation) is very popular technique to approximate stochastic integrals

of the form

$$\int e^{-Ml(\boldsymbol{\phi})} d\boldsymbol{\phi} \quad (6)$$

around a minimum $\hat{\boldsymbol{\phi}}$ of sufficiently smooth positive function l leading to

$$\int e^{-Ml(\boldsymbol{\phi})} d\boldsymbol{\phi} \rightarrow \left(\frac{2\pi}{M}\right)^{p/2} |H(\hat{\boldsymbol{\phi}})|^{-1/2} e^{-Ml(\hat{\boldsymbol{\phi}})} \quad (7)$$

as $M(\in \mathbb{R}^+) \rightarrow \infty$ (Wong, 2001). The evaluation of Hessian $H(\hat{\boldsymbol{\phi}})$ of l is needed. Second, when looking for a maximum-likelihood estimator $\hat{\boldsymbol{\phi}}^{MLE} := \underset{\boldsymbol{\phi}}{\operatorname{argmin}} l(\boldsymbol{\phi})$ one usually applies some optimization algorithm which requires many evaluations of gradient $\nabla l(\boldsymbol{\phi})$, such as Conjugate Gradient (CG) or Broyden-Fletcher-Goldfarb-Shanno (BFGS) (Bertsekas, 1999; Bazaraa et al., 2006).

Third, modern Monte Carlo Markov Chain (MCMC) samplers such as Metropolis-Adjusted Langevin algorithm (MALA) or Hamiltonian Monte Carlo (HMC) also require computing gradients or even Hessians of a log-likelihood with respect to the model parameters for every sample (Brooks et al., 2011).

For any of the above problems, the computation of the derivatives is a key operation and its speedup directly translates to the speedup of the whole algorithm. For example, in the case of the HMC sampler, the total speedup is roughly proportional to the speedup of the gradient computation of the log-likelihood.

1.2 Adjoint-state method

We will employ the *adjoint-state method* (ASM) to efficiently compute the first and the second derivative (gradient and Hessian) of (2) with respect to the model parameters $\boldsymbol{\phi}$.

The ASM is used in many different fields, such as control theory (Lions, 1971), data assimilation in meteorology (Lewis et al., 2006) or parameter identification (Melicher and Vrabel, 2013; Cimrak and Melicher, 2007). It is difficult to precisely trace its origin, since it is based on a general principle - the duality. Special dual problems or special test functions in general are used extensively in functional analysis for quite different tasks, e.g. in homogenization theory (Bensoussan et al., 1978). The idea of the ASM is to derive a special dual problem to the sensitivity equation of (1), which allows one to write the derivative(s) of (2) in a simple form which is inexpensive to evaluate. Usually, one obtains an inner product(s) in a suitable Hilbert space containing the dual state.

Even if the ASM is a classical method in many different fields, its applications in general statistical literature are rather scarce or could be even considered virtually non-existent. In our opinion, this is due to several reasons.

Mainly, it is a matter of need. Until recently the usual statistical models had only few parameters and the corresponding derivatives were easily evaluable using finite differences. On the other hand, the ASM was mainly used for problems, where derivatives with respect to infinite dimensional parameters are needed, such as the optimal

control of partial differential equations (PDEs) (Lions, 1971). For those problems, the ASM is the only viable way to compute the gradient of a cost functional.

The second possible reason is the lack of interdisciplinary publications in statistical literature with the fields where the ASM plays the role of a classical technique. An exception is the field of meteorology, where the relationship between differential models and stochastic processes is probably the most advanced (Lewis et al., 2006). Outside of this field, the worlds of the PDE-constrained inverse problems and Bayesian inference have been elegantly connected in (Martin et al., 2012). The paper presents a stochastic Newton method in which MCMC is sampling from a proposal density that builds a local Gaussian approximation based on local derivatives of the log posterior information. The authors exploit adjoint-based gradients and Hessians (as matrix-vector products). They argue that the effective dimension of a parameter estimation problem is often mesh-independent and consequently the Hessian can be approximated by a low-rank approximation computed using a Lanczos process.

The third and probably rather influential reason is that the results presented in literature regarding the ASM do not take into account the specifics of statistical estimation, particularly that the measurements can not be altered or interpolated in any way. In this paper, we present an ASM framework for ODE based statistical models, which recognizes and resolves this issue. The ODE case can be addressed in generality, which is not possible for PDE-constrained problems.

The dynamical models described by ODEs are rather widely used in science. They are simply indispensable for acquiring essential knowledge about complex biological systems (Murray, 2002; Draelants et al., 2012) as is the case for other fields studying intricate matters such as psychology and economics. In chemistry, regardless of the criticism (Gillespie, 1977), the reaction-rate equations¹ are still extensively employed. We are motivated by applications in PK/PD modeling and virology (Lavielle et al., 2011; Tornøe et al., 2004).

The sensitivity analysis of ODEs is well established in literature. Let us only mention a classical book on the optimal control of ODEs (Cesari, 1983). Moreover, many results that are intended for PDEs are directly applicable to ODEs, since from the mathematical point of view, an ODE could be simply seen as a PDE without a spatial differential operator. However, as already mentioned, the relevant results presented in literature, do not take into account the specifics of statistical estimation.

The ASM is usually applied in a PDE-constrained context. The fidelity between the data and the PDE-based model is measured in a Lebesgue space L^p -norm, particularly in L^2 sense, as is also the case of the above mentioned paper (Martin et al., 2012). It implies that the data are considered to be defined almost everywhere in the space or in the space-time in the case of time-dependent problems. This is however in a strong contrast with statistical philosophy. The measurements are ultimately discrete and sacred. By interpolating the measurements, new ones are generated and that can not be tolerated.

The main contribution of this paper is that it recognizes and resolves this problem. We show, that the discrete data y can be combined with the continuous model (1)

¹ Considering spatial phenomena such as diffusion and(or) convection leads to PDE models, see for example (Slodička and Balážová, 2010).

at the level of the likelihood functional (2). The resulting adjoint problem contains a Dirac delta source corresponding to individual measurement times. The developments are fully supported by rigorous proofs.

The subsequent numerical analysis shows that the ASM application for statistical estimation is far from obvious and more work is still needed to realize its deterministic setting efficiency. The main problem that remains to be solved is the inverse proportionality of ASM efficiency to the number of observations N . Our primary domain of interest is non-linear mixed effects modeling of population data (Lindstrom and Bates, 1990), where the number of measurement times N is usually very small but the number of model parameters p can be quite large, so the effect is rather subdued.

Since the ASM is usually applied in an infinite dimensional setting, as explained above, only results for first order derivatives are usually available. For the ODE case, we can supply the Hessian computation as well.

Another contribution is a highly efficient implementation of the results and its benchmarking with respect to finite differences and sensitivity equation approach.

We do not know about similar results in the literature.

Last but not least, this interdisciplinary paper aims to popularize this quite underused but potentially very useful method in the statistical community and help those working with ODE based models to compute the corresponding ODE-model sensitivities more efficiently. Recently, with the boom in general availability and dimensionality of data, ODE models with a high number of parameters are being employed. The evaluation of gradients becomes very costly and consequently various derivative-free methods have become more popular, see for example Delyon et al. (1999).

One of the domains, where the results presented on the next pages could be particularly appreciated is Systems biology. We refer the interested reader to Raue et al. (2013), where the quantitative dynamical modeling is assessed from a rather broad perspective. Notably relevant is a conclusion of the paper that multi-start deterministic parameter optimization using the sensitivity equations (see Section 2 here) for the calculation of derivatives significantly outperforms all other tested algorithms, including a number of stochastic optimization variants which do not make use of derivative information.

As we will show, for certain set-ups, the ASM significantly outperforms the sensitivity equation method for the computation of likelihood derivatives with respect to the parameters of the underlying ODE-model. The gradient can be computed with a cost (measured by model evaluations) that is essentially independent of the number of the parameters. The Hessian can be computed as well, with a linear rather than a quadratic cost in the number of parameters. Consequently, the use of ASM makes derivative-based (optimization) algorithms for certain setups even more competitive than presented in Raue et al. (2013).

The paper is structured as follows. In Section 2, we analyze the sensitivity of initial value problems (1) with respect to its parameter vector ϕ . In Section 3, we present in detail the approach to combine discrete data with a continuous model. Then in Section 4, we obtain the ASM for computing of the gradient and Hessian of (2) with respect to ϕ . Finally, the implementation is discussed in Section 5 and its efficiency is tested on a number of examples in Section 6.

2 Sensitivity of model

In this preparatory section we will discuss the well-posedness of the initial value problem (1) as well as the existence of its derivative with respect to the parameters ϕ . We follow the presentation in (Zeidler, 1985) with all the relevant notation, so we can be rather concise.

The first Gâteaux differential of a function f with respect to \mathbf{x} in direction \mathbf{h} is denoted by $\mathcal{D}f(\mathbf{x}; \mathbf{h})$. Then, let us denote by $\mathbf{s} := \mathcal{D}\mathbf{u}(\phi; \mathbf{h})$, i.e. the first Gâteaux differential (we will show it is Fréchet as well) of the model function \mathbf{u} with respect to the parameters ϕ in direction \mathbf{h} . If it exists, the formal differentiation of (1) yields that \mathbf{s} is the solution to the following initial value problem

$$\begin{aligned} d_t \mathbf{s} &= J_{\mathbf{u}}(\mathbf{f})\mathbf{s} + J_{\phi}(\mathbf{f})\mathbf{h}, \quad t \in [0, T], \\ \mathbf{s}(0) &= J_{\phi}(\mathbf{u}_0)\mathbf{h}, \end{aligned} \quad (8)$$

known as the *sensitivity equation*. Here $J_{\phi}(\mathbf{f}) : \mathbb{R}^p \rightarrow \mathbb{R}^m$ and $J_{\mathbf{u}}(\mathbf{f}) : \mathbb{R}^m \rightarrow \mathbb{R}^m$ are the Jacobians of the r.h.s. \mathbf{f} of the model (1) with respect to the model parameters ϕ and the state variables of the model \mathbf{u} , respectively. Similarly, $J_{\phi}(\mathbf{u}_0)$ denotes the Jacobian of the initial value with respect to the model parameters ϕ .

Let $\mathbf{e}_i, i = 1, \dots, p$ be the canonical basis in \mathbb{R}^p . Solving (8) for $\mathbf{h} = \mathbf{e}_i$ for each $i = 1, \dots, p$ yields $\mathbf{s} = (\frac{\partial \mathbf{u}_1}{\partial \phi_i}, \frac{\partial \mathbf{u}_2}{\partial \phi_i}, \dots, \frac{\partial \mathbf{u}_m}{\partial \phi_i})^t$, if the partial derivatives exist. It means that to compute the whole jacobian $J_{\phi}(\mathbf{u})$ one needs to integrate p initial value problems (8). The complexity is essentially identical to that of the first-order finite difference approximation, as will be confirmed in Section 6. The sensitivity equation approach is however still preferred if high accuracy is needed.

Let us restate the Theorem 4.D from (Zeidler, 1985) in our context.

Theorem 1 *Suppose that the mappings $\mathbf{f} : U \subseteq \mathbb{R} \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}^m$ and $\mathbf{u}_0 : V \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^m$ are $C^k, k \geq 1$ and that U, V are open sets containing $(0, \mathbf{u}_0(\phi_0), \phi_0)$ and ϕ_0 , respectively. Then:*

- (a) *There exists an interval $(-a, a), a > 0$, and an open neighborhood $U(\phi_0)$ such that the initial value problem (1) has exactly one solution for each $\phi \in U(\phi_0)$.*
- (b) *The mapping $(t, \phi) \mapsto \mathbf{u}(t; \phi)$ is C^k on $(-a, a) \times U(\phi_0)$, and (8) holds.*

Since our initial value problem (1) slightly differs from that of (Zeidler, 1985) and also for the completeness we present a proof in Appendix A.

From now on, any formal differentiation of \mathbf{u} with respect to the parameters ϕ is justified by Theorem 1. The theorem provides only a local result regarding the existence and the uniqueness of the solution \mathbf{u} of the ivp (1). Consequently, we have to assume that $T < a$.

3 Connecting the worlds

Measurements \mathbf{y}_i are acquired at discrete time points t_i . In statistics, these measurements should not be tampered with in any way and their interpolation would stand for augmentation.

On the other hand the model (1) is a continuous one and since the adjoint-state method (ASM) deals extensively with the model and the functional (2), it is necessary to work in continuous setting.

We will connect the discrete data and the continuous model on the level of the likelihood functional. One can write

$$\sum_i d(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi})) = \int_0^T \delta\{t - t_i\} d(\mathbf{y}(t), \mathbf{g}(t, \boldsymbol{\phi})) dt \quad (9)$$

where, by the classical misuse of notation, $\delta\{t - t_i\}$ is the Dirac delta function of the set of all measurement times t_i . In order to achieve that the above integral is well-defined, we will consider a small positive ε , such that the functions $\mathbf{y}(t) := \mathbf{y}_i, t \in (t_i - \varepsilon, t_i + \varepsilon)$ for all measurement times t_i are well defined. We emphasize, that by doing so, we do not generate new measurements. We merely assume an infinitesimally small interval of their validity. The $\mathbf{y}(t)$ -values outside of intervals $t \in (t_i - \varepsilon, t_i + \varepsilon)$ are irrelevant. For clarity, we extend the function $\mathbf{y}(t)$ outside of these intervals by linear interpolation to continuous functions on whole interval $[0, T]$ ².

For the well-posedness of (9), also the model $\mathbf{g}(t, \boldsymbol{\phi})$ has to be at least continuous around each t_i . Let us assume that $\mathcal{P} \in \mathcal{L}(\mathbb{R}^m, \mathbb{R}^n)$, i.e. \mathcal{P} is a linear operator from \mathbb{R}^m to \mathbb{R}^n . The linearity is a sufficient condition for the validity of

$$\mathcal{P}(\mathbf{u}(t, \boldsymbol{\phi})) - \mathcal{P}(\mathbf{u}_0(\boldsymbol{\phi})) = \int_0^t \mathcal{P}(\mathbf{f}(t, \boldsymbol{\phi}, \mathbf{u})) dt \quad (10)$$

which will be needed for the subsequent developments. From now on we write $\mathcal{P}\mathbf{u}$ instead of $\mathcal{P}(\mathbf{u})$. Let us point out that this assumption is usually not restrictive in practical applications. The eventual non-linear transformations can be applied a priori to the data \mathbf{y} or included in \mathbf{f} . Now, since the solution \mathbf{u} to (1) is at least continuously differentiable (Theorem 1) and a linear operator preserves continuity, $\mathbf{g}(t, \boldsymbol{\phi})$ is trivially continuous.

As a convenience, for any distribution d and sufficiently smooth function f , we denote by $\langle d, f \rangle$ the duality between them on the time interval $[0, T]$. We will also need the scalar product (\cdot, \cdot) in the Hilbert space $L^2([0, T])$. Using this notation, (9) can be rewritten as

$$\int_0^T \delta\{t - t_i\} d(\mathbf{y}(t), \mathbf{g}(t, \boldsymbol{\phi})) dt = \langle \delta\{t - t_i\}, d(\mathbf{y}(t), \mathbf{g}(t, \boldsymbol{\phi})) \rangle. \quad (11)$$

Let us introduce short notation $\{\delta\}$ for $\delta\{t - t_i\}$. At last, the equality (9) defines a seminorm on $C([0, T], \mathbb{R}^n)$, since the l.h.s. is a discrete norm. We denote this seminorm simply as $\|\cdot\|$.

Using the gluing notation above, we can prove the following lemma that allows us to evaluate the first differential of (2) using the solution \mathbf{s} to the sensitivity equation (8).

² Other continuous ‘‘interpolation’’ are possible such as piecewise-linear or by cubic splines, but they are less graphic.

Lemma 1 *Let the assumptions of Theorem 1 be fulfilled for $k = 1$ and let the metric d be C^1 . Then the functional (2) is Fréchet differentiable and the differential $\mathcal{D}l(\boldsymbol{\phi}; \mathbf{h})$ can be expressed as*

$$\mathcal{D}l(\boldsymbol{\phi}; \mathbf{h}) = \langle \{\delta\} d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})), \mathbf{s} \rangle, \quad (12)$$

where \mathbf{s} is the unique solution to sensitivity equation (8).

A proof can be found in Appendix A. Moreover, due to the linearity of (8), the differential $\mathcal{D}l(\boldsymbol{\phi}; \mathbf{h})$ can be easily written in the linearized form $\mathcal{D}l(\boldsymbol{\phi}; \mathbf{h}) = l'(\boldsymbol{\phi})\mathbf{h}$. Since we work in finite dimensional spaces, the expression

$$\nabla l(\boldsymbol{\phi}) \cdot \mathbf{h} := l'(\boldsymbol{\phi})\mathbf{h} \quad \text{for all } \mathbf{h} \in \mathbb{R}^p \quad (13)$$

well defines the gradient $\nabla l(\boldsymbol{\phi})$ of l as an element in \mathbb{R}^p for each fixed $\boldsymbol{\phi}$, i.e. $\nabla l : \mathbb{R}^p \rightarrow \mathbb{R}^p$.

As explained in Section 2, p initial value problems (8) have to be computed to evaluate $\nabla l(\boldsymbol{\phi})$ for some $\boldsymbol{\phi}$.

4 Adjoint-state method

In this Section we will introduce the *adjoint-state method* (ASM) for computation of the gradient and the Hessian of (2). The results are strongly influenced by the peculiar coupling between the discrete measurements and the continuous model (1). Let us directly present the main statement.

Theorem 2 *Let the assumptions of Lemma 1 be fulfilled. Then the first Fréchet differential in (12) can be also written as*

$$\mathcal{D}l(\boldsymbol{\phi}; \mathbf{h}) = -\mathbf{v}'(0)J_{\boldsymbol{\phi}}(\mathbf{u}_0)\mathbf{h} - (J_{\boldsymbol{\phi}}(\mathbf{f})\mathbf{h}, \mathbf{v}) \quad (14)$$

where \mathbf{v} is the unique solution to the following initial value problem

$$\begin{aligned} d_t \mathbf{v} &= -J_{\mathbf{u}}^t(\mathbf{f})\mathbf{v} + \{\delta\} d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})), \quad t \in [0, T], \\ \mathbf{v}(T) &= 0. \end{aligned} \quad (15)$$

A proof is again presented in Appendix A. Obviously, Equation (14) is written in linearized form. We get

$$\nabla l = -\mathbf{v}'(0)J_{\boldsymbol{\phi}}(\mathbf{u}_0) - (\mathbf{v}, J_{\boldsymbol{\phi}}(\mathbf{f})), \quad (16)$$

where the second term on the r.h.s. is a vectorial integral. This is a very efficient way to compute the gradient. One has to only integrate one adjoint problem (15) and evaluate the expression (16), i.e. to compute p scalar products in $L^2([0, T])$.

Let us discuss the result a little. The ivp (15) is a special ODE. First, it has absolutely no physical, chemical, biological or any other interpretation of the underlying scientific field of equation (1). The best way to look at it is that it is a special dual problem (see the proof) to the sensitivity equation (8), which allows us to efficiently compute the gradient of (2) (and the Hessian as well as we will see.) Then, it is a final

time problem to be integrated from T to the initial time 0. It is a linear ODE like the sensitivity equation. Its r.h.s. contains the term $\{\delta\}d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi}))$, which expresses how quickly the distance between the data and the model changes when changing the model variable \mathbf{u} .

Probably the most important fact to note about the ASM is that it operates at a higher level than the sensitivity equation method. It does not supply the derivative of the state \mathbf{u} , but directly the one of the likelihood functional (2). By considering the model together with (2), the efficiency can be achieved.

The numerical issues regarding the integration of (15) will be discussed in Section 5.

Example 1 Let us consider the distance (4) corresponding to the multivariate normal distribution of the residual errors. Then the derivative $d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi}))$ in the r.h.s of (15) reads

$$d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})) = -\mathcal{P}^t \Sigma^{-1}(\mathbf{y}_i - \mathbf{g}(t_i, \boldsymbol{\phi})). \quad (17)$$

The adjoint problem is dependent on the residual covariance matrix Σ and on the post-processing operator \mathcal{P} .

4.1 Evaluating Hessian

Let us depict the second Gateaux differential of a functional f with respect to \mathbf{x} in directions \mathbf{h}_1 and \mathbf{h}_2 as $\mathcal{D}^2 f(\mathbf{x}; \mathbf{h}_1, \mathbf{h}_2)$. Then, let us introduce notation $\boldsymbol{\zeta} := \mathcal{D}^2 \mathbf{u}(\boldsymbol{\phi}; \mathbf{h}_1, \mathbf{h}_2)$. We will show that $\boldsymbol{\zeta}$ is Fréchet as well. By formally differentiating (8) one more time with respect to $\boldsymbol{\phi}$ we obtain that $\boldsymbol{\zeta}$ is a solution to the following initial value problem

$$\begin{aligned} d_t \boldsymbol{\zeta} &= \mathbf{f}_{\phi\phi} \mathbf{h}_1 \mathbf{h}_2 + \mathbf{f}_{\phi\mathbf{u}} \mathbf{h}_1 \mathbf{s}_2 + \mathbf{f}_{\mathbf{u}\phi} \mathbf{s}_1 \mathbf{h}_2 \\ &\quad + \mathbf{f}_{\mathbf{u}\mathbf{u}} \mathbf{s}_1 \mathbf{s}_2 + J_{\mathbf{u}}(\mathbf{f}) \boldsymbol{\zeta}, \quad t \in [0, T], \\ \boldsymbol{\zeta}(0) &= (\mathbf{u}_0)_{\phi\phi} \mathbf{h}_1 \mathbf{h}_2 \end{aligned} \quad (18)$$

known as the *second sensitivity equation*. Here $\mathbf{s}_1, \mathbf{s}_2$ are the solutions to (8) for $\mathbf{h} = \mathbf{h}_1, \mathbf{h} = \mathbf{h}_2$, respectively. The second order derivatives in (18) are essentially three-dimensional tensors. In Appendix A the following lemma is proven.

Lemma 2 *Let the assumptions of Theorem 1 be fulfilled for $k = 2$ and let the metric d be C^2 . Then the second Fréchet differential of (2) with respect to the model parameters $\boldsymbol{\phi}$ can be written as*

$$\mathcal{D}^2 l(\boldsymbol{\phi}; \mathbf{h}_1, \mathbf{h}_2) = \langle \{\delta\}, d_{\mathbf{u}\mathbf{u}}^2 \mathbf{s}_1 \mathbf{s}_2 \rangle + \langle \boldsymbol{\zeta}, \{\delta\} d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})) \rangle, \quad (19)$$

where $\boldsymbol{\zeta}$ is the unique solution to (18). Moreover, the second term in (19) can be rewritten using the solution \mathbf{v} to (15) as

$$\begin{aligned} \langle \boldsymbol{\zeta}, \{\delta\} d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})) \rangle &= -(\mathbf{u}_0)_{\phi\phi} \mathbf{h}_1 \mathbf{h}_2 \cdot \mathbf{v}(0) - (\mathbf{f}_{\phi\phi} \mathbf{h}_1 \mathbf{h}_2, \mathbf{v}) \\ &\quad - (\mathbf{f}_{\phi\mathbf{u}} \mathbf{h}_1 \mathbf{s}_2, \mathbf{v}) - (\mathbf{f}_{\mathbf{u}\phi} \mathbf{s}_1 \mathbf{h}_2, \mathbf{v}) - (\mathbf{f}_{\mathbf{u}\mathbf{u}} \mathbf{s}_1 \mathbf{s}_2, \mathbf{v}). \end{aligned} \quad (20)$$

Solving (18) for $\mathbf{h}_1 = \mathbf{e}_i$, $\mathbf{h}_2 = \mathbf{e}_j$ for each $i = 1, \dots, p$, $j = 1, \dots, p$ yields $\boldsymbol{\zeta} = (\frac{\partial^2 \mathbf{u}_1}{\partial \boldsymbol{\phi}_i \partial \boldsymbol{\phi}_j}, \frac{\partial^2 \mathbf{u}_2}{\partial \boldsymbol{\phi}_i \partial \boldsymbol{\phi}_j}, \dots, \frac{\partial^2 \mathbf{u}_m}{\partial \boldsymbol{\phi}_i \partial \boldsymbol{\phi}_j})^t$. It means that to compute the Hessian $H_{\boldsymbol{\phi}}(\mathbf{u})$, one needs to integrate the second sensitivity equation (18) $p(p+1)/2$ times. For that one moreover needs to compute p sensitivities \mathbf{s}_i for each $\mathbf{h} = \mathbf{e}_i$, $i = 1, \dots, p$. The cost is essentially identical to that of the first order finite difference approximation. Again, it is beneficial if high accuracy is needed.

On the other hand, the evaluation of the Hessian $H_{\boldsymbol{\phi}}l$ of (2) via (20) requires us to only compute one adjoint problem (15), p sensitivity equations (8) and $p(p+1)/2$ times the four scalar products from (20). This is a very efficient and accurate way how to compute the Hessian.

As before with the gradient, we see that the ASM supplies the sensitivity at the level of the functional, not that of the underlying model state \mathbf{u} .

4.1.1 Hessian via adjoint with finite differences

Let us present an alternative way to efficiently compute the Hessian of (2), which is slightly less accurate than using (20) but much simpler to implement. The idea is to combine (14) with finite differences as follows

$$H_i(l(\boldsymbol{\phi})) \approx \frac{\nabla l(\boldsymbol{\phi} + h\mathbf{e}_i) - \nabla l(\boldsymbol{\phi})}{h}, \quad (21)$$

where H_i stands for the i -th column of H (or row) and h is a small positive real number. We recall that $\{\mathbf{e}_i : 1 \leq i \leq p\}$ is the canonical basis in \mathbb{R}^p . Each of p gradients $\nabla l(\boldsymbol{\phi} + h\mathbf{e}_i)$, $1 \leq i \leq p$ is computed using (14). Together $p+1$ adjoint initial value problems (15) need to be integrated.

5 Implementation of ASM

In this Section we will describe an implementation of the ASM presented in Section 4. At the core of the developments is the adjoint initial value problem (15). Even if it is a rather simple linear ODE-system, it is a quite difficult one to solve numerically because of its r.h.s. containing the Dirac delta function source term.

Our implementation closely follows the constructive proof of Theorem 2 in Appendix A. At each measurement time t_i , ODE-solver is stopped, $d_{\mathbf{u}}(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi}))$ is explicitly added to the current solution and then the integration is resumed (we solve a sequence of initial value problems (32) instead of the original ivp (15)).

Unfortunately, the repetitive restarting of the ODE solver has a negative impact on the performance. The derivative $d_{\mathbf{u}}(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi}))$ is added to the dynamical system at once via the initial condition $\mathbf{v}(t_i)$ and since the time derivative $d_t \mathbf{v}$ from (32) is proportional to $\mathbf{v}(t_i)$, it encounters a jump at each measurement point. Consequently, the steepness of the solution forces the ODE-solver to advance in many small time steps, which leads to a high number of iterations. We will see in Section 6, that the efficiency of ASM is indeed strongly dependent on the number of measurements.

However, equation (32) is a quite simple linear ODE-system, which should be exploitable in multiple ways. Although increasing the numerical efficiency of backward

integration while preserving the statistical rigor is a very interesting scientific goal, it is out of the scope of this contribution and left for the future research.

We tackle (15) using CVODES solver from the SUNDIALS package (Hindmarsh et al., 2005). CVODES is an extension of the CVODE code with both forward and backward sensitivity abilities (Serban and Hindmarsh, 2005).

The numerical experiments presented in Section 6 are computed in *DiffMEM* (Haber et al., 2016). It is a package for the fitting of mixed-effect models constrained by differential equations. The package is under active development by the authors and the algorithms presented in this paper are only a subset of its abilities.

At present, only ODE dynamical models are supported. *DiffMEM* employs the ODE solvers of CVODE for quick and robust solution of those models. It uses Eigen linear algebra package (library) to represent its internal memory containers and to solve underlying linear systems.

Remark 1 Let us imagine, we would not explicitly integrate the Dirac delta function out. It can be approximated in many different ways, but the most suitable from the statistical point of view (owing to the central limit theorem) is the approximation using Gaussian

$$\delta_i^\sigma(t) := \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(t-t_i)^2}{2\sigma^2}}, \quad (22)$$

where σ can be seen as a measure of the trust that the measurements have been taken precisely at the times t_i . Let us define $\delta^\sigma(t) := \sum_i \delta_i^\sigma(t)$. Using this approximation, the adjoint system (15) becomes

$$\begin{aligned} d_t \mathbf{v} &= -J'_u(\mathbf{f}) \mathbf{v} + \delta^\sigma(t) d_u(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})), \quad t \in [0, T], \\ \mathbf{v}(T) &= 0. \end{aligned} \quad (23)$$

Here, the adjoint problem (23) represents an interesting antagonism between the efficiency of the ASM and the statistical rigour one expects. The higher the trust in the measurement times, the smaller the σ and consequently higher the derivative of the r.h.s of (23) with respect to time which makes this dynamical system more and more difficult for an ODE solver to integrate.³

6 Numerical experiments

In this section we will consider for simplicity but without any loss of generality the Gaussian log-likelihood (5) with $\Sigma = I$.

For all the experiments, the ODE solvers' absolute accuracies are set to 10^{-14} and the relative ones to 10^{-10} . These accuracies are sufficient to remove considerations about accuracy of the ODE-solver from the analysis.

We will study the efficiency and robustness of the adjoint-state method (ASM) for computing the derivatives of the likelihood.

³ The variance σ^2 has here a purely ad hoc use for the above argument. We are not interested if it is prescribed or estimated from the data.

6.1 Linear model

To start, let us consider the classical linear ordinary differential equation (ODE)

$$\begin{aligned} d_t \mathbf{u} &= A\mathbf{u}, \quad t \in [0, T], \\ \mathbf{u}(0) &= \mathbf{u}_0, \end{aligned} \quad (24)$$

where A is a $d \times d$ -dimensional matrix, the elements of which represent the model parameters. This simple model is ideal toy-example to comprehend the importance of ASM for models with high dimensional parametric space.

Let us consider diagonal matrix A . Then the dimensions of the parametric space and of the solution coincide ($p = m$). Moreover, we can easily calculate the exact solution

$$\mathbf{u}_i = \mathbf{u}_{0,i} e^{\phi_i t}, \quad i = 1, \dots, p, \quad (25)$$

where $\phi_i = A_{ii}$.

We consider 13 different dimensions of ϕ , ranging from 2 to 122. For each of them we have randomly generated 100 parameter-samples ϕ as follows

$$\phi_i \sim U[-1.1, -0.1], \quad 1 \leq i \leq p. \quad (26)$$

We set $\mathbf{u}_0 = \mathbf{1}$. The number of observation time points N , regularly spread in $[0, 100]$, is set constant to 11.

The corresponding synthetic data \mathbf{y}_ϕ are also perturbed as follows:

$$\mathbf{y}_i = \mathbf{y}_{\phi,i} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 0.01[\max(\mathbf{y}_\phi)]^2 I), \quad 1 \leq i \leq N, \quad (27)$$

where $\max(\mathbf{y}_\phi) \in \mathbb{R}$ is taken through all the data and $I \in \mathbb{R}^{n \times n}$ is the identity matrix. We would like to emphasize that any reasonable perturbation leads to the same results. It is only important that the data are perturbed outside of the log-likelihood mode.

We have computed the gradient of likelihood using finite differences (FD), the ASM approach (14) and using the sensitivity equation (SE) (8). We recall that the last two approaches are implemented using the CVODES forward- and backward-sensitivity abilities and all the common settings are identical to make comparison as sound as possible. The results are presented in Figure 1.

The timings of FD, ASM and of SE are presented in Figure 1(a). The first conclusion is that our implementation of SE-approach is optimal since the timings of FD and SE more or less coincide. Actually SE is always a slightly faster method. Given the significantly higher accuracy of SE with respect to FD (1(d)), it renders FD-approach redundant.

Somewhere around 10 parameters ASM becomes on average more efficient than SE. Moreover, given the non-parametric prediction intervals based on the 100 samples, it is from around 15 parameters virtually always more time-efficient than SE. This reasoning is conservative since it does not take the correlation between ASM and SE timings into account. Moreover, the variance of timings is for $\dim(\phi) > 10$ lower for ASM than for SE, see Figure 1(b), rendering timing predictions for ASM more reliable.

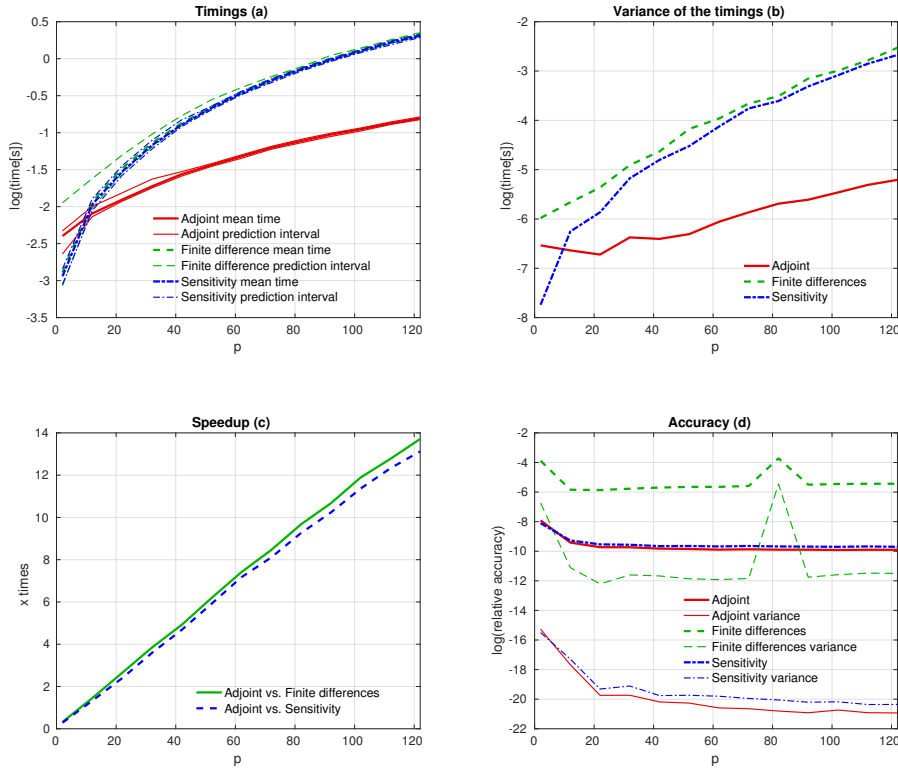


Fig. 1 Linear ODE model: comparison of different methods for log-likelihood gradient computation with respect to the growing parameter space dimension of the ODE model

The time efficiency of both ASM and SE is negatively impacted by exclusive use of dense matrices in *DiffMEM*. The equations (8), (15), (14) require evaluation of Jacobians $J_{\mathbf{u}}(\mathbf{f})$ and $J_{\phi}(\mathbf{f})$. These are extremely sparse⁴. More importantly, because of the dense matrix implementation only rather small systems can be solved. Sparse matrix implementation is planned for the future versions of *DiffMEM*. Both ASM and SE are influenced to the same degree and the relative comparison holds.

The speedup of ASM vs. SE (1(c)) is roughly linear in the number of the parameters but it slows down slightly for higher parameter dimensions. The suspected cause here is the cost of memory access when CVODES evaluates the forward solution \mathbf{u} during the backwards integration of (15).

The accuracy of both ASM and SE with respect to the exact gradient of the likelihood (5) is presented in Figure 1(d). Both methods achieve virtually identical accuracy since the forward- and backward- solvers use the same relative and absolute tolerances.

Now we will examine the efficiency of ASM and SE with respect to the number of time observations. We fix the dimension of the problem at e.g. $\dim(\phi) = 50$. The

⁴ The diagonal system (24) is the most sparse system one can think of.

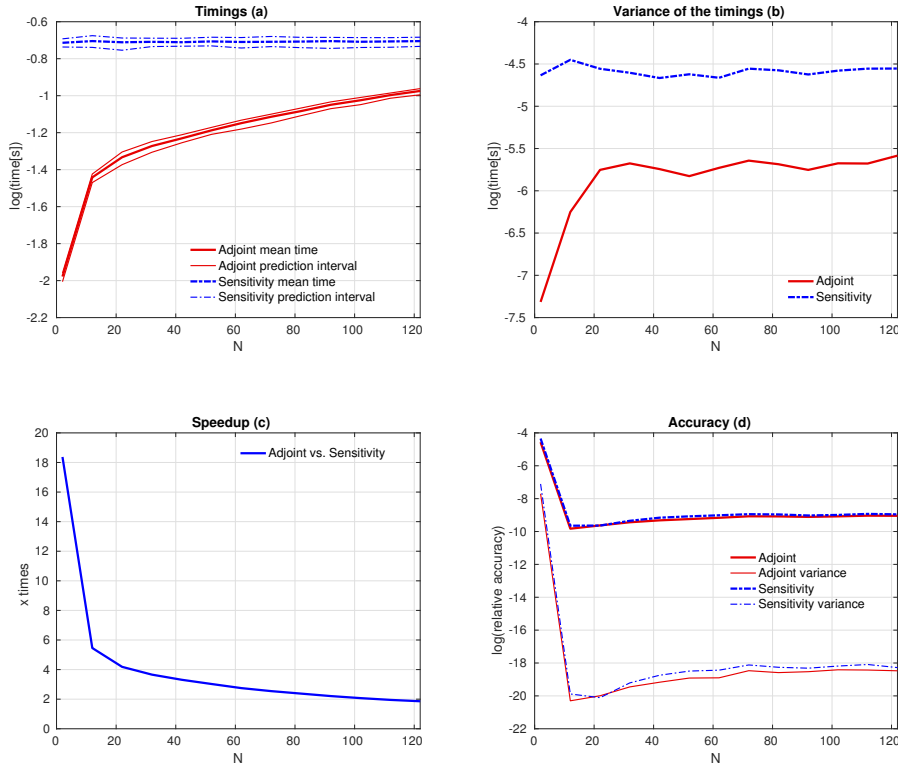


Fig. 2 Linear ODE model: comparison of different methods for log-likelihood gradient computation with respect to the growing number of the number of time observations

number of time observations $\dim(y)$ in $[0, 100]$ fluctuates between 2 and 122 in 12 steps. For each number we again compute 100 gradients using (8) and (14). The parameters ϕ are again generated using (26). The results are presented in Figure 2.

The SE approach efficiency is essentially independent of the number of time observations (2(a)). The ASM efficiency however decreases with increasing number of observations. As explained in Section 5, the backward adjoint integrator needs after each data point a large number of small time steps to account for the steepness of the adjoint solution \mathbf{v} . The negative impact is the most clearly visible in Figure 2(c). For many practically relevant problems (PK/PD, virology), the number of measurements is rather low, making this issue less pronounced. Anyhow, increasing the numerical efficiency of backward integration while preserving the statistical rigor is obviously a very interesting direction for future research.

Implicitly, since for the diagonal linear model $p = n$, also the dimension of the solution space plays a role. But we do not compare the speed and accuracy of the different methods with respect to m or n , since their complexities with respect to these are the same.

Now, again using the problem (24), we will illustrate the efficiency of computing the Hessian of (5) with respect to ϕ employing the expression (20). We compare this

(SA approach) with Hessian evaluated using the finite difference approximation (FD) and the one computed using (21) (FA).

Remark 2 First-order Gauss-Newton approximation of the Hessian, where the second term in (19) is neglected, is not included in the comparison. When the model does not yet well approximate the data, the second order term (20) can be arbitrarily large with respect to the first order term in (19). This is a well known fact but often overlooked. Let us return to the linear model (24) for a deeper insight. In this case, the first order approximation F of the Hessian of the likelihood H is

$$\begin{aligned} F_{k,k} &= -\sum_{i=1}^N e^{\phi_k t_i} e^{\phi_k t_i} t_i^2, \quad k = 1, \dots, p \\ F_{k,l} &= 0, \quad k \neq l \end{aligned} \quad (28)$$

and the second order term S is

$$\begin{aligned} S_{k,k} &= -\sum_i^N \left(e^{\phi_k t_i} - \mathbf{y}_i \right) e^{\phi_k t_i} t_i^2, \quad k = 1, \dots, p \\ S_{k,l} &= 0, \quad k \neq l. \end{aligned} \quad (29)$$

We see that the first order approximation F carries absolutely no information about how far the solution is from the data. The Gauss-Newton approximation error can thus be arbitrarily large when \mathbf{y} is not well approximated by the solution $e^{\phi t}$, especially for the values corresponding to small measurement times.

This is for example exploited in the well-known Levenberg-Marquardt method for the least-square minimization (Moré, 1978), which dynamically switches from the gradient descent method (GD) to the Gauss-Newton (GN) method. The GD is used to get sufficiently close to a minimum, so that the GN approximation is reliable.

The overall setup stays identical to the one used for the gradient, i.e. the one corresponding to Figure 1. For convenience, we consider a shorter parameter range $\dim(\phi) \in [2, 52]$, since the finite difference approximation of Hessian, to which we compare the ASM, has quadratic complexity in p . It makes the experiments more time prohibiting in comparison to the gradient. The results are depicted in Figure 3.

FD-adjoint, i.e. approximating Hessians using (21), is the most time-efficient approach (Figure 3(a)). The speedup with respect to the finite difference approximation (FD) is linear in $\dim(\phi)$ as expected (Figure 3(c)). The FD-adjoint Hessian accuracy is usually sufficient (Figure 3(d)) and moreover it is behaving well as a function of the dimension of the parametric space.

If higher accuracy up to machine precision is desirable, one can compute Hessian using (20), i.e. using SA-approach. Our SA-implementation is however clearly slower than the FD-adjoint despite of a rather optimal coding. FD-adjoint is superior time-wise mainly due to its simplicity.

To conserve space, we do not include any experiments for Hessian with respect to the number of measurement times N . But obviously, for the Hessian computed via FD-adjoint, the results for gradient are directly applicable. We will analyze the

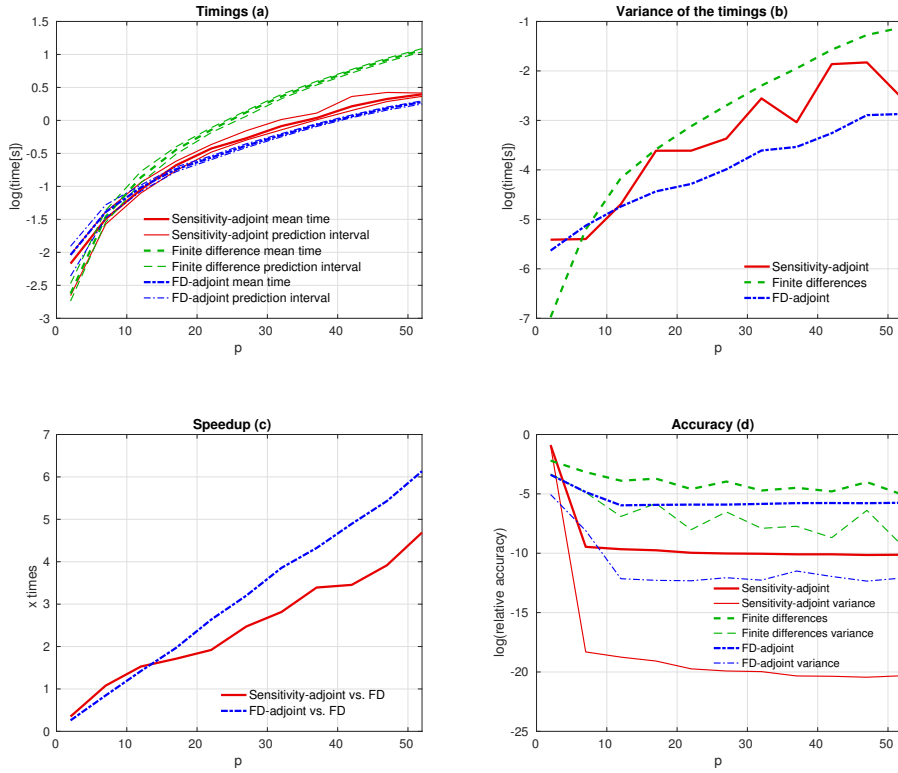


Fig. 3 Linear ODE model: comparison of different methods for log-likelihood Hessian computation with respect to the growing parameter space dimension of the ODE model

dependence on N for the following model in Section 6.2. Here we have focused on p -scaling only, which cannot be tested for the realistic model.

6.2 Latent dynamic HIV model

Now we are going to assess the efficiency and accuracy of ASM for a more realistic model - latent dynamic HIV model from Lavielle et al. (2011):

$$\begin{aligned}
 d_t T_{NI} &= \lambda - (1 - \eta_{NRTI})\gamma T_{NI} V_I - \mu_{NI} T_{NI}, \\
 d_t T_L &= (1 - \pi)(1 - \eta_{NRTI})\gamma T_{NI} V_I - \alpha_L T_L - \mu_L T_L, \\
 d_t T_A &= \pi(1 - \eta_{NRTI})\gamma T_{NI} V_I + \alpha_L T_L - \mu_A T_A, \\
 d_t V_I &= (1 - \eta_{PI})\rho T_A - \mu_V V_I, \\
 d_t V_{NI} &= \eta_{PI}\rho T_A - \mu_V V_{NI},
 \end{aligned} \tag{30}$$

where T_{NI} is the number of not-infected CD4 cells, T_L of latent infected CD4 cells and T_A the number of active infected CD4 cells producing new virions. The number of infectious viruses is V_I and the non-infectious V_{NI} . The 11 parameters ϕ represent

Table 1 Parameters of the latent dynamic HIV model

ϕ	λ	γ	μ_{NI}	μ_L	μ_A	μ_V	p	α_L	π	η_{NRTI}	η_{PI}
ϕ_0	2.61	.0021	.0085	.0092	.289	30	641	1.6×10^{-5}	.443	.90	.99

mostly rates of change. The two of them $\eta_{NRTI}, \eta_{PI} \in [0, 1]$ represent the efficacies of two types of antiviral therapies. The available measurements \mathbf{y}_i are restricted to the cumulative counts of CD4 cells and the viremia, i.e. $V_{ij} = V_I + V_{NI}$ and $T_{ij} = T_{NI} + T_L + T_A$ respectively. For the details see Lavielle et al. (2011). We have $p = 11$, $m = 5$ and $n = 2$.

The setup of the experiments stays rather similar to the previous ones. The parameters are generated randomly around ϕ_0 which is presented in Table 1 as follows:

$$\phi_i \sim U[0.95\phi_{0,i}, 1.05\phi_{0,i}], \quad 1 \leq i \leq p. \quad (31)$$

The efficacies η_{NRTI} and η_{PI} can be sometimes generated out of the allowed range $[0, 1)$. We project them back:

$$\phi_i = \min(\phi_i, 0.999), \quad i \in \{10, 11\}.$$

The corresponding synthetic measurements \mathbf{y} are perturbed using (27).

For the HIV model p is fixed and we can supply the results only with respect to N . We again observe in Figure 4(a) that the efficiency of (16) is strongly dependent on the number of observations. For up to 5 observations it is more efficient than the sensitivity equation (SE) approach. Thus even for models with a relatively small number of ODE parameters, the ASM approach for the computation of the gradient of (2) can be advised for certain applications, such as mixed effects modeling in pharmacokinetics and pharmacodynamics, as in (Lavielle et al., 2011). However, for models with a few parameters and a high number of observation points, the SE approach is clearly more efficient. Accuracy-wise, both approaches are equivalent (Figure 4(d)).

In Figure 5 the corresponding results for the Hessian computation of (2) are presented. Three ways are compared: finite difference (FD) approximation, the ASM approach (SA) using (19) and (20) and the mixed approach (FA) using (21).

First, again as for the diagonal linear model in Section 6.1, the efficiency of the FD approximation is virtually independent of the number of measurements (Figure 5(a)). This is not the case for the SA and FA approaches. However, due to its simplicity, the mixed FA approach is clearly more efficient than SA. It is more efficient than the FD approach up to 5 measurements, which corresponds to the previous results for the computation of the gradient.

The accuracies in Figure 5(d) are compared to the results of SA approach, as no exact solution is available. For the linear model (24), this approach was shown to be accurate up to the machine precision. The mixed FA approach achieves stable accuracies around 10^{-6} , two orders of magnitude better than the full finite difference approximation.

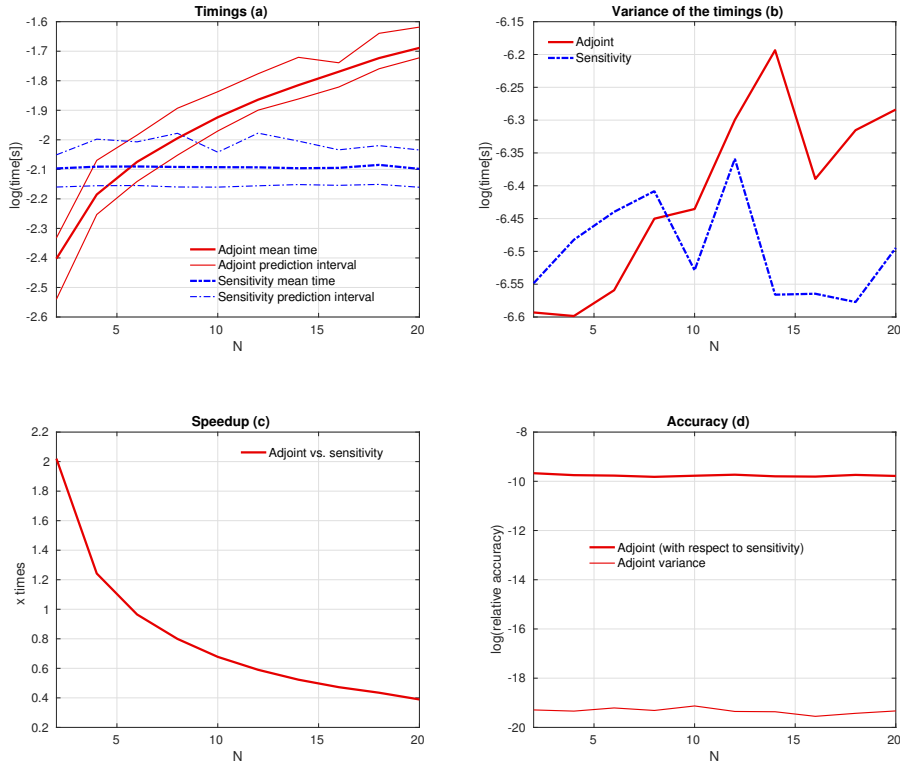


Fig. 4 HIV model: comparison of different methods for log-likelihood gradient computation with respect to the growing number of time observations

7 Conclusions

We have derived and analyzed the adjoint-state method for computation of the gradient and the Hessian of likelihood functionals for time series data modeled by ordinary differential equations. We interfaced the discrete data and the continuous model on the level of likelihood functional, using the concept of point-wise distributions. The resulting adjoint problem (15) then contains a Dirac delta source corresponding to individual measurement times. The developments are fully supported by the corresponding theoretical results. The implementation of a solver to (15) closely follows the constructive proof of its well-posedness.

Then, we compared the efficiency of the resulting ASM with finite differences and sensitivity equation (SE) approaches, both for the gradient and the Hessian. First, the implementation of SE approach is so efficient, that it renders the finite difference approximation practically obsolete, due to its superior accuracy. Second, the ASM efficiency is dependent on the number of measurement times, which is not the case for SE approach. For models with a high-number of parameters and a small number of measurement times, the ASM is a clear winner. It starts to be competitive even

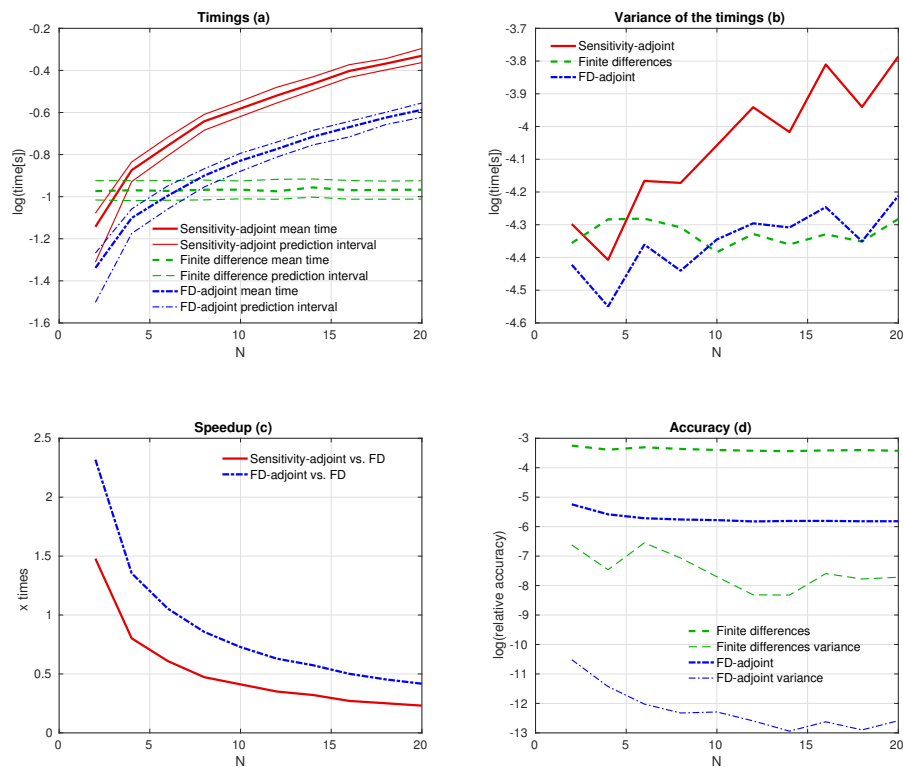


Fig. 5 HIV model: comparison of different methods for log-likelihood Hessian computation with respect to the growing number of time observations

for rather small models like the latent dynamic HIV model from Section 6.2 (11 parameters, 6 measurement times).

In future, we plan a sparse matrix code rewrite, which would allow for solution to bigger ODE systems and also a computationally more efficient implementation. Then, the preconditioning of Newton solver step during the CVODES integration of (15) is an interesting possibility to speed up the ASM (Knoll and Keyes, 2004).

Acknowledgements We would like to thank Xavier Woot de Trixhe from Janssen Pharmaceutica for numerous very interesting discussions on PK/PD, virology, biological pathways modeling, NLMEMs and on life in general. They were an important source of motivation and provided a view from a different perspective. And we would like to thank the anonymous reviewers for their substantial input, enhancing the quality of the paper.

8 Supplemental materials

Two external files are provided:

1. A document titled: “Supplemental material A: help with reproducing of the results presented in Fast derivatives of likelihood functionals for ODE based models using adjoint-state method”.
2. A simple and extensively commented R-implementation of the ASM. It computes the gradient of the Gaussian log-likelihood (5) with respect to the parameters of a pharmacokinetic two-compartment model. We set $\Omega = I$. The gradient is computed by the sensitivity equation method as well and those can be compared. But not in the terms of time efficiency. It is not a simple task to efficiently implement the ASM since the adjoint equation (15) is integrated backward in time and the solver has to compute its r.h.s. dependent on the forward solution. This has to be cached by the solver. The R-solution uses a simple linear scheme which is far from optimal but illustrative. As described in Section 5, the *DiffMEM* implementation uses heavily the capabilities of the SUNDIALS package.
The R-code is independently understandable but nevertheless references the relevant formulas of this paper.

A Proofs

Proof of Theorem 1

First, when compared with Theorem 4.D from Zeidler (1985) we work with $X = \mathbb{R}^m$ and $P = \mathbb{R}^p$. Those are complete normed vector spaces i.e. they are Banach spaces. Second, the initial condition is dependent only on parameter ϕ , not on any free parameter y as in Theorem 4.D.

Set $J = [-1, 1]$. Let us do the following rescaling: $t = sa$, $\mathbf{z}(s) := \mathbf{u}(as) - \mathbf{u}_0(\phi)$ for all $s \in J$. Then (1) is equivalent to

$$\mathbf{z}'(s) - a\mathbf{f}(as, \mathbf{z}(s) + \mathbf{u}_0(\phi), \phi) = \mathbf{0} \quad \text{for all } s \in J, \mathbf{z}(0) = \mathbf{0}.$$

This can be written as an operator equation $F(\mathbf{z}, a, \phi) = 0$ with the operator $F : \mathbf{Z} \times \mathbf{A} \rightarrow \mathbf{W}$ and spaces $\mathbf{Z} = \{\mathbf{z} \in C^1(J, \mathbb{R}^m) : \mathbf{z}(0) = \mathbf{0}\}$, $\mathbf{W} = C(J, \mathbb{R}^m)$. The space \mathbf{A} contains all the parameters (a, ϕ) , i.e. $\mathbf{A} = \mathbb{R} \times \mathbb{R}^p$.

Set $\mathbf{q} = (0, 0, \phi)$. Both F and F_z are continuous at \mathbf{q} . Obviously, $F(\mathbf{q}) = \mathbf{0}$ and $F_z(\mathbf{q})\mathbf{z} = \mathbf{z}'$. The crucial observation is that for every $\mathbf{w} \in \mathbf{W}$, there exists exactly one $\mathbf{z} \in \mathbf{Z}$ with $\mathbf{z}' = \mathbf{w}$, namely $\mathbf{z}(s) = \int_0^s \mathbf{w}(t) dt$. Hence $F_z(\mathbf{q}) : \mathbf{Z} \rightarrow \mathbf{W}$ is bijective. The implicit function theorem yields the conclusions (see e.g. Theorem 4.B in Zeidler (1985)). ■

Proof of Lemma 1

After realizing that l depends on ϕ only through \mathbf{u} , (12) is formally a direct application of the chain rule ($d_{\mathbf{u}}$ denotes the derivative of the metric with respect to the model state \mathbf{u} .)

The r.h.s. of (12) is a well-posed finite expression. First, we have assumed that the metric is sufficiently smooth, thus $d_{\mathbf{u}}$ is continuous. Second, \mathbf{s} is continuous as well owing to Theorem 1. A distribution can be rescaled by any at least continuous function, as here $\{\tilde{\delta}\}$ by $d_{\mathbf{u}}$.

Thus, the first differential on the l.h.s of (12) exists as well. It is moreover continuous, i.e. it is Fréchet.
5 ■

⁵ An alternative argumentation could employ equivalence (9).

Proof of Theorem 2

First, let us assume that we have already constructed a unique solution \mathbf{v} to (15) up to a certain measurement point t_i . The adjoint problem is solved backwards in time. Consequently, we will construct its prolongation on $[t_i, t_{i-1})$.

Let \mathbf{v}_i^+ be the ODE solution just before integrating the measurement at time t_i , i.e. at time t_i^+ . We simply stop the integration at t_i^+ , add $d_{\mathbf{u}}(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi}))$ to \mathbf{v}_i^+ and solve

$$\begin{aligned} d_t \mathbf{v} &= -J_{\mathbf{u}}^t(\mathbf{f}) \mathbf{v}, \quad t \in (t_i, t_{i-1}), \\ \mathbf{v}(t_i) &= \mathbf{v}_i^+ + d_{\mathbf{u}}(\mathbf{y}_i, \mathbf{g}(t_i, \boldsymbol{\phi})). \end{aligned} \quad (32)$$

This is a simple linear ODE with a continuous coefficient $J_{\mathbf{u}}^t(\mathbf{f})$, since $f \in \mathbf{C}^1$. The classical results yield the global solution on (t_i, t_{i-1}) (see e.g. Theorem 5.1 and Theorem 5.2 from (Coddington and Levinson, 1955).) This concludes the proof of the existence and uniqueness.

Now, we prove (14). Let us without a loss of generality assume that there are no measurements in times 0 and T . It is a well-known result of theory of distributions (in the sense of functional analysis), that the classical integration by part formula

$$\int_0^T d_t \mathbf{v} \mathbf{w} dt = [\mathbf{v} \mathbf{w}]_0^T - \int_0^T \mathbf{v} d_t \mathbf{w} dt \quad (33)$$

is valid for $\mathbf{w} \in \mathbf{C}^1$ even if the derivative $d_t \mathbf{v}$ exists on $[0, T]$ only in a *weak* sense, i.e. almost everywhere. Actually, (33) is the definition of the weak derivative of \mathbf{v} taking only $\mathbf{w} \in \mathbf{C}_0^1([0, T])$. Consequently, since $\mathbf{s} \in \mathbf{C}^1([0, T])$, we can safely proceed as follows

$$\begin{aligned} \langle \mathbf{s}, \{\delta\} d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})) \rangle &\stackrel{(15)}{=} \langle \mathbf{s}, d_t \mathbf{v} + J_{\mathbf{u}}^t(\mathbf{f}) \mathbf{v} \rangle \\ &\stackrel{(8)}{=} -\mathbf{v}'(0) J_{\boldsymbol{\phi}}(\mathbf{u}_0) \mathbf{h} - (d_t \mathbf{s} - J_{\mathbf{u}}(\mathbf{f}) \mathbf{s}, \mathbf{v}) \\ &\stackrel{(8)}{=} -\mathbf{v}'(0) J_{\boldsymbol{\phi}}(\mathbf{u}_0) \mathbf{h} - (J_{\boldsymbol{\phi}}(\mathbf{f}) \mathbf{h}, \mathbf{v}). \blacksquare \end{aligned} \quad (34)$$

Proof of Lemma 2

The existence and uniqueness of $\boldsymbol{\zeta}$ is a direct results of Theorem 1. Now, (20) is derived as follows:

$$\begin{aligned} \langle \boldsymbol{\zeta}, \{\delta\} d_{\mathbf{u}}(\mathbf{y}, \mathbf{g}(t, \boldsymbol{\phi})) \rangle &\stackrel{(15)}{=} \langle \boldsymbol{\zeta}, d_t \mathbf{v} + J_{\mathbf{u}}^t(\mathbf{f}) \mathbf{v} \rangle \\ &\stackrel{(15), (18)}{=} -(\mathbf{u}_0)_{\boldsymbol{\phi} \boldsymbol{\phi}} \mathbf{h}_1 \mathbf{h}_2 \cdot \mathbf{v}(0) \\ &\quad - (d_t \boldsymbol{\zeta}, \mathbf{v}) + (J_{\mathbf{u}}(\mathbf{f}) \boldsymbol{\zeta}, \mathbf{v}). \end{aligned} \quad (35)$$

This after substituting for $J_{\mathbf{u}}(\mathbf{f}) \boldsymbol{\zeta}$ from (18) directly yields (20). Analogically to the proof of Theorem 2, we needed $\boldsymbol{\zeta} \in \mathbf{C}^1([0, T])$ to be able to integrate by parts. \blacksquare

References

- Bazaraa MS, Sherali HD, Shetty CM (2006) Nonlinear programming: theory and algorithms, 3rd edn. John Wiley & Sons, Hoboken
- Bensoussan A, Lions J, Papanicolaou G (1978) Asymptotic analysis for periodic structures. North-Holland Pub. Co, Amsterdam
- Bertsekas DP (1999) Nonlinear programming. Athena scientific, Belmont
- Brooks S, Gelman A, Jones G, Meng XL (2011) Handbook of Markov Chain Monte Carlo. CRC Press, Boca Raton
- Cesari L (1983) Optimization—theory and applications: problems with ordinary differential equations, Applications of Mathematics, vol 17. Springer-Verlag, New York

- Cimrák I, Melicher V (2007) Sensitivity analysis framework for micromagnetism with application to the optimal shape design of magnetic random access memories. *Inverse Problems* 23(2):563–588
- Coddington EA, Levinson N (1955) *Theory of ordinary differential equations*. Tata McGraw-Hill Education, New York
- Delyon B, Lavielle M, Moulines E (1999) Convergence of a stochastic approximation version of the EM algorithm. *The Annals of Statistics* 27(1):94–128
- Draelants D, Broeckhove J, Beemster GTS, Vanroose W (2012) Numerical bifurcation analysis of the pattern formation in a cell based auxin transport model. *Journal of Mathematical Biology* 67(5):1279–1305
- Gillespie DT (1977) Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry* 81(25):2340–2361
- Haber T, Melicher V, Michiels N, Kovac T, Nemeth B, Claes J (2016) DiffMEM, <https://bitbucket.org/tomhaber/diffmem/branch/analysis>
- Hindmarsh AC, Brown PN, Grant KE, Lee SL, Serban R, Shumaker DE, Woodward CS (2005) SUNDIALS: Suite of Nonlinear and Differential/Algebraic Equation Solvers. *ACM Trans Math Softw* 31(3):363–396
- Knoll D, Keyes D (2004) Jacobian-free Newton–Krylov methods: a survey of approaches and applications. *Journal of Computational Physics* 193(2):357–397
- Lavielle M, Samson A, Karina Fermin A, Mentré F (2011) Maximum Likelihood Estimation of Long-Term HIV Dynamic Models and Antiviral Response. *Biometrics* 67(1):250–259
- Lewis JM, Lakshmiarahan S, Dhall S (2006) Dynamic data assimilation: a least squares approach, *Encyclopedia of Mathematics and its Applications*, vol 13. Cambridge University Press, Cambridge
- Lindstrom MJ, Bates DM (1990) Nonlinear mixed effects models for repeated measures data. *Biometrics* 46(3):pp. 673–687
- Lions JL (1971) *Optimal control of systems governed by partial differential equations*. Springer, Berlin
- Martin J, Wilcox L, Burstedde C, Ghattas O (2012) A Stochastic Newton MCMC Method for Large-Scale Statistical Inverse Problems with Application to Seismic Inversion. *SIAM Journal on Scientific Computing* 34(3):A1460–A1487
- Melicher V, Vrabel' V (2013) On a continuation approach in Tikhonov regularization and its application in piecewise-constant parameter identification. *Inverse Problems* 29(11):115,008
- Moré JJ (1978) The Levenberg-Marquardt algorithm: Implementation and theory. In: Watson G (ed) *Numerical Analysis*, Lecture Notes in Mathematics, vol 630, Springer, Berlin, pp 105–116
- Murray JD (2002) *Mathematical Biology I: An Introduction*, Interdisciplinary Applied Mathematics, vol 17, 3rd edn. Springer-Verlag, New York
- Raue A, Schilling M, Bachmann J, Matteson A, Schelke M, Kaschek D, Hug S, Kreutz C, Harms BD, Theis FJ, Klingmüller U, Timmer J (2013) Lessons learned from quantitative dynamical modeling in systems biology. *PLOS ONE* 8(9):1–17
- Serban R, Hindmarsh AC (2005) CVODES: the sensitivity-enabled ODE solver in SUNDIALS. In: ASME 2005 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference, American Society of Mechanical Engineers, pp 257–269
- Slodička M, Balážová A (2010) Decomposition method for solving multi-species reactive transport problems coupled with first-order kinetics applicable to a chain with identical reaction rates. *Journal of Computational and Applied Mathematics* 234(4):1069–1077, proceedings of the Thirteenth International Congress on Computational and Applied Mathematics (ICCAM-2008), Ghent, Belgium, 7–11 July, 2008
- Tornøe CW, Agersø, H, Jonsson E, Madsen H, Nielsen HA (2004) Non-linear mixed-effects pharmacokinetic/pharmacodynamic modelling in NLME using differential equations. *Computer Methods and Programs in Biomedicine* 76(1):3–40
- Wong R (2001) Asymptotic approximations of integrals, *Classics in applied mathematics*, vol 34. SIAM, Boston
- Zeidler E (1985) *Nonlinear Functional Analysis and Its Applications: Fixed point theorems*. *Nonlinear Functional Analysis and Its Applications*, Springer-Verlag, New York