



**HAL**  
open science

# f-SAEM: A fast Stochastic Approximation of the EM algorithm for nonlinear mixed effects models

Belhal Karimi, Marc Lavielle, Éric Moulines

► **To cite this version:**

Belhal Karimi, Marc Lavielle, Éric Moulines. f-SAEM: A fast Stochastic Approximation of the EM algorithm for nonlinear mixed effects models. Computational Statistics and Data Analysis, inPress, 10.1016/j.csda.2019.07.001 . hal-01958248

**HAL Id: hal-01958248**

**<https://inria.hal.science/hal-01958248>**

Submitted on 17 Dec 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# f-SAEM: A fast Stochastic Approximation of the EM algorithm for nonlinear mixed effects models

Belhal Karimi<sup>a,b,\*</sup>, Marc Lavielle<sup>a,b</sup>, Eric Moulines<sup>a,b</sup>

<sup>a</sup>*CMAP, Ecole Polytechnique, route de Saclay, 91120 Palaiseau, France*

<sup>b</sup>*INRIA Saclay, 1 Rue Honoré d'Estienne d'Orves, 91120 Palaiseau, France*

---

## Abstract

The ability to generate samples of the random effects from their conditional distributions is fundamental for inference in mixed effects models. Random walk Metropolis is widely used to perform such sampling, but this method is known to converge slowly for high dimensional problems, or when the joint structure of the distributions to sample is spatially heterogeneous. We propose an independent Metropolis–Hastings (MH) algorithm based on a multidimensional Gaussian proposal that takes into account the joint conditional distribution of the random effects and does not require any tuning. Indeed, this distribution is automatically obtained thanks to a Laplace approximation of the incomplete data model. We show that such approximation is equivalent to linearizing the structural model in the case of continuous data. Numerical experiments based on simulated and real data illustrate the performance of the proposed methods. In particular, we show that the suggested MH algorithm can be efficiently combined with a stochastic approximation version of the EM algorithm for maximum likelihood estimation in nonlinear mixed effects models.

*Keywords:* MCMC, Stochastic approximation, EM, mixed effects, Laplace approximation

---

## 1. Introduction

Mixed effects models are often adopted to take into account the inter-individual variability within a population (see (Lavielle, 2014) and the references therein). Consider a study with  $N$  individuals from a same population. The vector of observations  $y_i$  associated to each individual  $i$  is assumed to be a realisation of a random variable which

---

\*Corresponding author

*Email addresses:* `belhal.karimi@inria.fr` (Belhal Karimi), `marc.lavielle@inria.fr` (Marc Lavielle), `eric.moulines@polytechnique.edu` (Eric Moulines)

depends on a vector of random individual parameters  $\psi_i$ . Then, inference on the individual parameter  $\psi_i$  amounts to estimate its conditional distribution given the observed data  $y_i$ .

When the model is a linear (mixed effects) Gaussian model, then this conditional distribution is a normal distribution that can explicitly be computed (Verbeke, 1997). For more complex distributions and models, Monte Carlo methods must be used to approximate this conditional distribution. Most often, direct sampling from this conditional distribution is inefficient and it is necessary to resort to a Markov chain Monte Carlo (MCMC) method for obtaining random samples from this distribution.

Note that generating random samples from  $p_i(\psi_i|y_i)$  is useful for several tasks to avoid approximation of the model, such as linearisation or Laplace method. Such tasks include the estimation of the population parameters  $\theta$  of the model, either by a maximum likelihood approach, i.e. by maximizing the observed incomplete data likelihood  $p(y_1, \dots, y_N; \theta)$  using the Stochastic Approximation of the EM algorithm (SAEM) algorithm combined with a MCMC procedure (Kuhn and Lavielle, 2004), or by a Bayesian method, i.e. by estimating  $p(\theta|y_1, \dots, y_N)$ . Lastly, sampling from the conditional distributions  $p_i(\psi_i|y_i)$  is also known to be useful for model building. Indeed, in (Lavielle and Ribba, 2016), the authors argue that methods for model assessment and model validation, whether graphical or based on statistical tests, must use samples of the conditional distribution  $p_i(\psi_i|y_i)$  to avoid bias.

Designing a fast mixing sampler for these distributions is therefore of utmost importance. The most common MCMC method for nonlinear mixed effects (NLME) models is the *random walk Metropolis* (RWM) algorithm (Robert and Casella, 2010; Roberts et al., 1997; Lavielle, 2014). This method is implemented in software tools such as Monolix, NONMEM, the saemix R package (Comets et al., 2017) and the nlmeftsa Matlab function.

Despite its simplicity, it has been successfully used in many classical examples of pharmacometrics. Nevertheless, it can show its limitations when the parameter space to explore becomes large or when the dependency structure of the individual parameters is complex. In particular, maintaining an optimal acceptance rate (advocated in (Roberts and Rosenthal, 1997)) most often implies very small moves and therefore a very large number of iterations. Therefore, if we want to adapt the MCMC to high-dimensional probability distributions of practical interest, we need to better use the geometry of the target distribution.

The Metropolis-adjusted Langevin algorithm (MALA) uses evaluations of the gradient of the target density for proposing new states which are accepted or rejected using the Metropolis–Hastings algorithm (Roberts and Tweedie, 1996; Stramer and Tweedie, 1999). Several variations have been proposed for improving the behaviour of MALA by incorporating more information about the properties of the target distribution in the proposal, see for instance (Girolami and Calderhead, 2011; Allasonniere and Kuhn, 2013; Durmus et al., 2017).

Hamiltonian Monte Carlo (HMC) is another MCMC algorithm that exploits informa-

tion about the geometry of the target distribution in order to efficiently explore the space by selecting transitions that can follow contours of high probability mass (Betancourt, 2017). The No-U-Turn Sampler (NUTS) is an extension to HMC that allows an automatic and optimal selection of some of the settings required by the algorithm, (Brooks et al., 2011; Hoffman and Gelman, 2014).

Nevertheless, these methods may be difficult to use in practice, and are computationally involved, in particular when the structural model is a complex ODE based model. The algorithm we propose is an independent Metropolis-Hastings (IMH) algorithm, but for which the proposal is a Gaussian approximation of the target distribution. For general data model (i.e. categorical, count or time-to-event data models or continuous data models), the Laplace approximation of the incomplete pdf  $\mathbf{p}_i(y_i)$  leads to a Gaussian approximation of the conditional distribution  $\mathbf{p}_i(\psi_i|y_i)$ .

In the special case of continuous data, linearisation of the model leads, by definition, to a Gaussian linear model for which the conditional distribution of the individual parameter  $\psi_i$  given the data  $y_i$  is a multidimensional normal distribution that can be computed. Therefore, we design an independent sampler using this multivariate Gaussian distribution to sample from target conditional distribution.

The paper is organised as follows. Mixed effects models for continuous and non-continuous data are presented in Section 2. The standard MH for NLME models is described in Section 3. The proposed method, called the nlme-IMH, is introduced in Section 4. The f-SAEM, a combination of this new method with the SAEM algorithm for estimating the population parameters of the model is specified in Section 5. Numerical examples illustrate, in Section 6, the practical performances of the proposed method, both on a continuous pharmacokinetics (PK) model and a time-to-event example. A Monte Carlo study confirms that this new SAEM algorithm shows a faster convergence to the maximum likelihood estimate.

## 2. Mixed Effect Models

### 2.1. Population approach and hierarchical models

In the sequel, we adopt a population approach, where we consider  $N$  individuals and  $n_i$  observations per individual  $i$ . The set of observed data is  $y = (y_i, 1 \leq i \leq N)$  where  $y_i = (y_{ij}, 1 \leq j \leq n_i)$  are the observations for individual  $i$ . For the sake of clarity, we assume that each observation  $y_{ij}$  takes its values in some subset of  $\mathbb{R}$ . The distribution of the  $n_i$ -vector of observations  $y_i$  depends on a vector of individual parameters  $\psi_i$  that takes its values in a subset of  $\mathbb{R}^p$ .

We assume that the pairs  $(y_i, \psi_i)$  are mutually independent and consider a parametric framework: the joint distribution of  $(y_i, \psi_i)$  is denoted by  $\mathbf{p}_i(y_i, \psi_i; \theta)$ , where  $\theta$  is the vector of parameters of the model. A natural decomposition of this joint distribution reads

$$\mathbf{p}_i(y_i, \psi_i; \theta) = \mathbf{p}_i(y_i|\psi_i; \theta)\mathbf{p}_i(\psi_i; \theta), \quad (1)$$

where  $\mathbf{p}_i(y_i|\psi_i; \theta)$  is the conditional distribution of the observations given the individual parameters, and where  $\mathbf{p}_i(\psi_i; \theta)$  is the so-called population distribution used to describe the distribution of the individual parameters within the population.

A particular case of this general framework consists in describing each individual parameter  $\psi_i$  as the sum of a typical value  $\psi_{\text{pop}}$  and a vector of individual random effects  $\eta_i$ :

$$\psi_i = \psi_{\text{pop}} + \eta_i . \quad (2)$$

In the sequel, we assume that the random effects are distributed according to a multivariate Gaussian distribution:  $\eta_i \sim_{\text{i.i.d.}} \mathcal{N}(0, \Omega)$ .

Several extensions of model (2) are also possible. We can assume for instance that the transformed individual parameters are normally distributed:

$$u(\psi_i) = u(\psi_{\text{pop}}) + \eta_i , \quad (3)$$

where  $u$  is a strictly monotonic transformation applied on the individual parameters  $\psi_i$ . Examples of such transformation are the logarithmic function (in which case the components of  $\psi_i$  are log-normally distributed), the logit and the probit transformations (Lavielle, 2014). In the following, we either use the original parameter  $\psi_i$  or the Gaussian transformed parameter  $u(\psi_i)$ .

Another extension of model (2) consists in introducing individual covariates in order to explain part of the inter-individual variability:

$$u(\psi_i) = u(\psi_{\text{pop}}) + C_i \beta + \eta_i , \quad (4)$$

where  $C_i$  is a matrix of individual covariates. Here, the fixed effects are the vector of coefficients  $\beta$  and the vector of typical parameters  $\psi_{\text{pop}}$ .

## 2.2. Continuous data models

A regression model is used to express the link between continuous observations and individual parameters:

$$y_{ij} = f_i(t_{ij}, \psi_i) + \varepsilon_{ij} , \quad (5)$$

where  $y_{ij}$  is the  $j$ -th observation for individual  $i$  measured at time  $t_{ij}$ ,  $\varepsilon_{ij}$  is the residual error. It is assumed that for all time  $t$ ,  $\psi \rightarrow f(t, \psi)$  is twice differentiable in  $\psi$ .

We start by assuming that the residual errors are independent and normally distributed with zero-mean and a constant variance  $\sigma^2$ . Let  $t_i = (t_{ij}, 1 \leq n_i)$  be the vector of observation times for individual  $i$ . Then, the model for the observations reads:

$$y_i | \psi_i \sim \mathcal{N}(f_i(\psi_i), \sigma^2 \mathbf{Id}_{n_i \times n_i}) ,$$

where

$$f_i(\psi_i) = (f_i(t_{i,1}, \psi_i), \dots, f_i(t_{i,n_i}, \psi_i)) . \quad (6)$$

If we assume that  $\psi_i \sim_{\text{i.i.d.}} \mathcal{N}(\psi_{\text{pop}}, \Omega)$ , then the parameters of the model are  $\theta = (\psi_{\text{pop}}, \Omega, \sigma^2)$ .

An extension of this model consists in assuming that the variance of the residual errors is not constant over time:

$$\varepsilon_{ij} \sim \mathcal{N}(0, g(t_{ij}, \psi_i)^2). \quad (7)$$

Such extension includes proportional error models ( $g = bf$ ) and combined error models ( $g = a + bf$ ) (Lavielle, 2014) but the proposed method remains the same whatever the residual error model is.

### 2.3. Noncontinuous data models

Noncontinuous data models include categorical data models (Savic et al., 2011; Agresti, 1990), time-to-event data models (Mbogning et al., 2015; Andersen, 2006), or count data models (Savic et al., 2011).

A categorical outcome  $y_{ij}$  takes its value in a set  $\{1, \dots, L\}$  of  $L$  categories. Then, the model is defined by the conditional probabilities ( $\mathbb{P}(y_{ij} = \ell | \psi_i), 1 \leq \ell \leq L$ ), that depend on the vector of individual parameters  $\psi_i$  and may be a function of the time  $t_{ij}$ .

In a time-to-event data model, the observations are the times at which events occur. An event may be one-off (e.g., death, hardware failure) or repeated (e.g., epileptic seizures, mechanical incidents). To begin with, we consider a model for a one-off event. The survival function  $S(t)$  gives the probability that the event happens after time  $t$ :

$$S(t) \triangleq \mathbb{P}(T > t) = \exp \left\{ - \int_0^t h(u) du \right\}, \quad (8)$$

where  $h$  is called the hazard function. In a population approach, we consider a parametric and individual hazard function  $h(\cdot, \psi_i)$ .

The random variable representing the time-to-event for individual  $i$  is typically written  $T_i$  and may possibly be right-censored. Then, the observation  $y_i$  for individual  $i$  is

$$y_i = \begin{cases} T_i & \text{if } T_i \leq \tau_c \\ "T_i > \tau_c" & \text{otherwise,} \end{cases} \quad (9)$$

where  $\tau_c$  is the censoring time and " $T_i > \tau_c$ " is the information that the event occurred after the censoring time.

For repeated event models, times when events occur for individual  $i$  are random times ( $T_{ij}, 1 \leq j \leq n_i$ ) for which conditional survival functions can be defined:

$$\mathbb{P}(T_{ij} > t | T_{i(j-1)} = t_{i(j-1)}) = \exp \left\{ - \int_{t_{i(j-1)}}^t h(u, \psi_i) du \right\}. \quad (10)$$

Here,  $t_{ij}$  is the observed value of the random time  $T_{ij}$ . If the last event is right censored, then the last observation  $y_{i,n_i}$  for individual  $i$  is the information that the censoring time

has been reached " $T_{i,n_i} > \tau_c$ ". The conditional pdf of  $y_i = (y_{ij}, 1 \leq n_i)$  reads (see (Lavielle, 2014) for more details)

$$p_i(y_i|\psi_i) = \exp \left\{ - \int_0^{\tau_c} h(u, \psi_i) du \right\} \prod_{j=1}^{n_i-1} h(t_{ij}, \psi_i). \quad (11)$$

### 3. Sampling from conditional distributions

#### 3.1. The conditional distribution of the individual parameters

Once the conditional distribution of the observations  $p_i(y_i|\psi_i; \theta)$  and the marginal distribution of the individual parameters  $\psi_i$  are defined, the joint distribution  $p_i(y_i, \psi_i; \theta)$  and the conditional distribution  $p_i(\psi_i|y_i; \theta)$  are implicitly specified. This conditional distribution  $p_i(\psi_i|y_i; \theta)$  plays a crucial role for inference in NLME models.

One of the main task is to compute the maximum likelihood (ML) estimate of  $\theta$

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta, y), \quad (12)$$

where  $\mathcal{L}(\theta, y) = \log p(y; \theta)$ . In NLME models, this optimization is solved by using a surrogate function defined as the conditional expectation of the complete data log-likelihood (McLachlan and Krishnan, 2007). The SAEM is an iterative procedure for ML estimation that requires to generate one or several samples from this conditional distribution at each iteration of the algorithm.

Once the ML estimate  $\hat{\theta}_{\text{ML}}$  has been computed, the observed Fisher information matrix

$$I(\hat{\theta}_{\text{ML}}, y) = -\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}_{\text{ML}}, y) \quad (13)$$

can be derived thanks to the Louis formula (Louis, 1982) which expresses  $I(\hat{\theta}_{\text{ML}}, y)$  in terms of the conditional expectation and covariance of the complete data log-likelihood. Such procedure also requires to sample from the conditional distributions  $p_i(\psi_i|y_i; \hat{\theta}_{\text{ML}})$  for all  $i \in \llbracket 1, N \rrbracket$ .

Samples from the conditional distributions might also be used to define several statistical tests and diagnostic plots for models assessment. It is advocated in (Lavielle and Ribba, 2016) that such samples should be preferred to the modes of these distributions (also called *Empirical Bayes Estimate*(EBE), or *Maximum a Posteriori Estimate*), in order to provide unbiased tests and plots. For instance, a strong bias can be observed when the EBEs are used for testing the distribution of the parameters or the correlation between random effects.

In short, being able to sample individual parameters from their conditional distribution is essential in nonlinear mixed models. It is therefore necessary to design an efficient method to sample from this distribution.

### 3.2. The Metropolis-Hastings Algorithm

Metropolis-Hasting (MH) algorithm is a powerful MCMC procedure widely used for sampling from a complex distribution (Brooks et al., 2011). To simplify the notations, we remove the dependency on  $\theta$ . For a given individual  $i \in \llbracket 1, N \rrbracket$ , the MH algorithm, to sample from the conditional distribution  $\mathbf{p}_i(\psi_i|y_i)$ , is described as:

---

#### Algorithm 1 Metropolis-Hastings algorithm

---

**Initialization:** Initialize the chain sampling  $\psi_i^{(0)}$  from some initial distribution  $\xi_i$ .

**Iteration k:** given the current state of the chain  $\psi_i^{(k-1)}$ :

1. Sample a candidate  $\psi_i^c$  from a proposal distribution  $q_i(\cdot|\psi_i^{(k-1)})$ .
2. Compute the MH ratio:

$$\alpha(\psi_i^{(k-1)}, \psi_i^c) = \frac{\mathbf{p}_i(\psi_i^c|y_i)}{\mathbf{p}_i(\psi_i^{(k-1)}|y_i)} \frac{q_i(\psi_i^{(k-1)}|\psi_i^c)}{q_i(\psi_i^c|\psi_i^{(k-1)})}. \quad (14)$$

3. Set  $\psi_i^{(k)} = \psi_i^c$  with probability  $\min(1, \alpha(\psi_i^{(k-1)}, \psi_i^c))$  (otherwise, keep  $\psi_i^{(k)} = \psi_i^{(k-1)}$ ).
- 

Under weak conditions,  $(\psi_i^{(k)}, k \geq 0)$  is an ergodic Markov chain whose distribution converges to the target  $\mathbf{p}_i(\psi_i|y_i)$  (Brooks et al., 2011).

Current implementations of the SAEM algorithm in Monolix (Chan et al., 2011), saemix (R package) (Comets et al., 2017), nlmeftsa (Matlab) and NONMEM (Beal and Sheiner, 1980) mainly use the same combination of proposals. The first proposal is an independent MH algorithm which consists in sampling the candidate state directly from the prior distribution of the individual parameter  $\psi_i$ . The MH ratio then reduces to  $\mathbf{p}_i(y_i|\psi_i^c)/\mathbf{p}_i(y_i|\psi_i^{(k)})$  for this proposal.

The other proposals are component-wise and block-wise random walk procedures (Metropolis et al., 1953) that updates different components of  $\psi_i$  using univariate and multivariate Gaussian proposal distributions. These proposals are centered at the current state with a diagonal variance-covariance matrix; the variance terms are adaptively adjusted at each iteration in order to reach some target acceptance rate (Atchadé and Rosenthal, 2005; Lavielle, 2014).

Nevertheless, those proposals have several drawbacks: such procedure is not suitable for sampling distributions in high dimension; and it fails to take into account the nonlinear dependence structure of the individual parameters.

A way to alleviate these problems is to use a proposal distribution derived from a discretised Langevin diffusion whose drift term is the gradient of the logarithm of the target density leading to the Metropolis Adjusted Langevin Algorithm (MALA) (Roberts



and Tweedie, 1996; Stramer and Tweedie, 1999). The MALA proposal is given by:

$$\psi_i^c \sim \mathcal{N}(\psi_i^{(k)} - \gamma \nabla_{\psi_i} \log \mathbf{p}_i(\psi_i^{(k)} | y_i), 2\gamma), \quad (15)$$

where  $\gamma$  is a positive stepsize. These methods appear to scale well in high dimension but still do not take into consideration the multidimensional structure of the individual parameters. Recent works include efforts in that direction, such as the Anisotropic MALA for which the covariance matrix of the proposal depends on the gradient of the target measure (Allasonniere and Kuhn, 2013), the Tamed Unadjusted Langevin Algorithm (Brosse et al., 2017) based on the coordinate-wise taming of superlinear drift coefficients and a multidimensional extension of the Adaptive Metropolis algorithm (Haario et al., 2001) simultaneously estimating the covariance of the target measure and coercing the acceptance rate, see (Vihola, 2012).

The MALA algorithm is a special instance of the Hybrid Monte Carlo (HMC), introduced in (Neal et al., 2011); see (Brooks et al., 2011) and the references therein, and consists in augmenting the state space with an auxiliary variable  $p$ , known as the velocity in Hamiltonian dynamics. This algorithm belongs to the class of data augmentation methods. Indeed, the potential energy is augmented with a kinetic energy, function of an added auxiliary variable. The MCMC procedure then consists in sampling from this augmented posterior distribution. All those methods aim at finding the proposal  $q$  that accelerates the convergence of the chain. Unfortunately they are computationally involved (in high dimension the computation of the gradient or the Hessian can be overwhelming) and can be difficult to implement (stepsizes and numerical derivatives need to be tuned and implemented).

We see in the next section how to define a multivariate Gaussian proposal for both continuous and noncontinuous data models, that is easy to implement and that takes into account the multidimensional structure of the individual parameters in order to accelerate the MCMC procedure.

#### 4. A multivariate Gaussian proposal

In this section, we assume that the individual parameters  $(\psi_1, \dots, \psi_N)$  are independent and normally distributed with mean  $(m_1, \dots, m_N)$  and covariance  $\Omega$ . The MAP estimate, for individual  $i$ , is the value of  $\psi_i$  that maximizes the conditional distribution  $\mathbf{p}_i(\psi_i | y_i, \theta)$ :

$$\hat{\psi}_i = \arg \max_{\psi_i \in \mathbb{R}^p} \mathbf{p}_i(\psi_i | y_i) = \arg \max_{\psi_i \in \mathbb{R}^p} \mathbf{p}_i(y_i | \psi_i) \mathbf{p}_i(\psi_i) \quad (16)$$

##### 4.1. Proposal based on Laplace approximation

For both continuous and noncontinuous data models, the goal is to find a simple proposal, a multivariate Gaussian distribution in our case, that approximates the target distribution  $\mathbf{p}_i(\psi_i | y_i)$ . It is well known, see (Mengersen and Tweedie, 1996; Roberts and Rosenthal, 2011) that the Independent Sampler is geometrically ergodic if and only if,

for a given  $\epsilon$ ,  $\inf_{\psi \in \mathbb{R}^p} q(\psi_i)/\mathbf{p}_i(\psi_i|y_i) \geq \epsilon > 0$  where  $q(\psi_i)$  is the proposal distribution. More generally, it is shown in (Roberts and Rosenthal, 2011) that the mixing rate in total variation depends on the expectation of the acceptance ratio under the proposal distribution which is also directly related to the ratio of proposal to the target. This observation naturally suggests to find a proposal which approximated the target. (de Freitas et al., 2001) advocates the use a multivariate Gaussian distribution whose parameters are obtained by minimizing the Kullback-Leibler divergence between a multivariate Gaussian variational candidate distribution and the target distribution. In (Andrieu and Thoms, 2008) and the references therein, an adaptative Metropolis algorithm is studied and reconciled to a KL divergence minimisation problem where the resulting multivariate Gaussian distribution can be used as a proposal in a IMH algorithm. Authors note that although this proposal might be a sensible choice when it approximates well the target, it can fail when the parametric form of the proposal is not sufficiently rich. Thus, other parametric forms can be considered and it is suggested in (Andrieu et al., 2006) to consider mixtures, finite or infinite, of distributions belonging to the exponential family.

In general, this optimization step is difficult and computationally expensive since it requires to approximate (using Monte Carlo integration for instance) the integral of the log-likelihood with respect to the variational candidate distribution.

**Independent proposal 1.** *We suggest a Laplace approximation of this conditional distribution as described in (Rue et al., 2009) which is the multivariate Gaussian distribution with mean  $\hat{\psi}_i$  and variance-covariance*

$$\Gamma_i = \left( -\mathbf{H}_{\hat{\psi}_i} + \Omega^{-1} \right)^{-1} \quad (17)$$

where  $\mathbf{H}_{\hat{\psi}_i} \in \mathbb{R}^{p \times p}$  is the Hessian of  $\log(\mathbf{p}_i(y_i|\psi_i))$  evaluated at  $\hat{\psi}_i$ .

Mathematical details for computing this proposal are postponed to Appendix A. We use this multivariate Gaussian distribution as a proposal in our IMH algorithm introduced in the next section, for both continuous and noncontinuous data models.

We shall now see another method to derive a Gaussian proposal distribution in the specific case of continuous data models (see (5)).

#### 4.2. Nonlinear continuous data models

When the model is described by (5), the approximation of the target distribution can be done twofold: either by using the Laplace approximation, as explained above, or by linearizing the structural model  $f_i$  for any individual  $i$  of the population. using (5) and (16), the MAP estimate can thus be derived as:

$$\hat{\psi}_i = \arg \min_{\psi_i \in \mathbb{R}^p} \left( \frac{1}{\sigma^2} \|y_i - f_i(\psi_i)\|^2 + (\psi_i - m_i)' \Omega^{-1} (\psi_i - m_i) \right). \quad (18)$$

where  $f_i(\psi_i)$  is defined by (6) and  $A'$  is the transpose of the matrix  $A$ .

We linearize the structural model  $f_i$  around the MAP estimate  $\hat{\psi}_i$ :

$$f_i(\psi_i) \approx f_i(\hat{\psi}_i) + \mathbf{J}_{f_i(\hat{\psi}_i)}(\psi_i - \hat{\psi}_i), \quad (19)$$

where  $\mathbf{J}_{f_i(\hat{\psi}_i)} \in \mathbb{R}^{n_i \times p}$  is the Jacobian of  $f_i$  evaluated at  $\hat{\psi}_i$ . Defining  $z_i := y_i - f_i(\hat{\psi}_i) + \mathbf{J}_{f_i(\hat{\psi}_i)} \hat{\psi}_i$ , this expansion yields the following linear model:

$$z_i = \mathbf{J}_{f_i(\hat{\psi}_i)} \psi_i + \varepsilon_i. \quad (20)$$

We can directly use the definition of the conditional distribution under a linear model (see (53) in Appendix B) to get an expression of the conditional covariance  $\Gamma_i$  of  $\psi_i$  given  $z_i$  under (20):

$$\Gamma_i = \left( \frac{\mathbf{J}'_{f_i(\hat{\psi}_i)} \mathbf{J}_{f_i(\hat{\psi}_i)}}{\sigma^2} + \Omega^{-1} \right)^{-1}. \quad (21)$$

Using (18),  $\hat{\psi}_i$  satisfies:

$$-\frac{\mathbf{J}'_{f_i(\hat{\psi}_i)}}{\sigma^2} (y_i - f_i(\hat{\psi}_i)) + \Omega^{-1}(\hat{\psi}_i - m_i) = 0, \quad (22)$$

which leads to the definition of the conditional mean  $\mu_i$  of  $\psi_i$  given  $z_i$ , under the linearized model, by:

$$\mu_i = \Gamma_i \frac{\mathbf{J}'_{f_i(\hat{\psi}_i)}}{\sigma^2} (y_i - f_i(\hat{\psi}_i) + \mathbf{J}_{f_i(\hat{\psi}_i)} \hat{\psi}_i + \Omega^{-1} m_i) \quad (23)$$

$$= \Gamma_i \left( \Omega^{-1}(\hat{\psi}_i - m_i) + \frac{\mathbf{J}'_{f_i(\hat{\psi}_i)} \mathbf{J}_{f_i(\hat{\psi}_i)}}{\sigma^2} \hat{\psi}_i + \Omega^{-1} m_i \right) = \Gamma_i \Gamma_i^{-1} \hat{\psi}_i = \hat{\psi}_i. \quad (24)$$

We note that the mode of the conditional distribution of  $\psi_i$  in the nonlinear model (5) is also the mode and the mean of the conditional distribution of  $\psi_i$  in the linear model (20).

**Independent proposal 2.** *In the case of continuous data models, we propose to use the multivariate Gaussian distribution, with mean  $\hat{\psi}_i$  and variance-covariance matrix  $\Gamma_i$  defined by (21) as a proposal for an independent MH algorithm avoiding the computation of an Hessian matrix.*

We can note that linearizing the structural model is equivalent to using the Laplace approximation with the expected information matrix. Indeed:

$$\mathbb{E}_{y_i|\hat{\psi}_i} \left( -\mathbf{H}_{l(\hat{\psi}_i)} \right) = \frac{\mathbf{J}'_{f_i(\hat{\psi}_i)} \mathbf{J}_{f_i(\hat{\psi}_i)}}{\sigma^2}. \quad (25)$$

**Remarks:**

1. When the model is linear, the probability of accepting a candidate generated with this proposal is equal to 1.
2. If we consider a more general error model,  $\varepsilon_i \sim \mathcal{N}(0, \Sigma(t_i, \psi_i))$  that may depend on the individual parameters  $\psi_i$  and the observation times  $t_i$ , then the conditional variance-covariance matrix reads:

$$\Gamma_i = \left( \mathbf{J}'_{f_i(\hat{\psi}_i)} \Sigma(t_i, \hat{\psi}_i)^{-1} \mathbf{J}_{f_i(\hat{\psi}_i)} + \Omega^{-1} \right)^{-1}. \quad (26)$$

3. In the model (3), the transformed variable  $\phi_i = u(\psi_i)$  follows a normal distribution. Then a candidate  $\phi_i^c$  is drawn from the multivariate Gaussian proposal with parameters:

$$\mu_i = \hat{\phi}_i, \quad (27)$$

$$\Gamma_i = \left( \frac{\mathbf{J}'_{f_i(u^{-1}(\hat{\phi}_i))} \mathbf{J}_{f_i(u^{-1}(\hat{\phi}_i))}}{\sigma^2} + \Omega^{-1} \right)^{-1}, \quad (28)$$

where  $\hat{\phi}_i = \arg \max_{\phi_i \in \mathbb{R}^p} \mathbf{p}_i(\phi_i | y_i)$  and finally the candidate vector of individual parameters is set to  $\psi_i^c = u^{-1}(\phi_i^c)$

These approximations of the conditional distribution  $\mathbf{p}_i(\psi_i | y_i)$  lead to our nlme-IMH algorithm, an Independent Metropolis-Hastings (IMH) algorithm for NLME models. For all individuals  $i \in \llbracket 1, N \rrbracket$ , the algorithm is defined as:

---

**Algorithm 2** The nlme-IMH algorithm

---

**Initialization:** Initialize the chain sampling  $\psi_i^{(0)}$  from some initial distribution  $\xi_i$ .

**Iteration t:** Given the current state of the chain  $\psi_i^{(t-1)}$ :

1. Compute the MAP estimate:

$$\hat{\psi}_i^{(t)} = \arg \max_{\psi_i \in \mathbb{R}^p} \mathbf{p}_i(\psi_i | y_i). \quad (29)$$

2. Compute the covariance matrix  $\Gamma_i^{(t)}$  using either (17) or (21).
3. Sample a candidate  $\psi_i^c$  from a the independent proposal  $\mathcal{N}(\hat{\psi}_i^{(t)}, \Gamma_i^{(t)})$  denoted  $q_i(\cdot | \hat{\psi}_i^{(t)})$ .
4. Compute the MH ratio:

$$\alpha(\psi_i^{(t-1)}, \psi_i^c) = \frac{\mathbf{p}_i(\psi_i^c | y_i) q_i(\hat{\psi}_i^{(t)} | \psi_i^c)}{\mathbf{p}_i(\psi_i^{(t-1)} | y_i) q_i(\psi_i^c | \hat{\psi}_i^{(t)})}. \quad (30)$$

5. Set  $\psi_i^{(t)} = \psi_i^c$  with probability  $\min(1, \alpha(\psi_i^c, \psi_i^{(t-1)}))$  (otherwise, keep  $\psi_i^{(t)} = \psi_i^{(t-1)}$ ).
-

This method shares some similarities with (Titsias and Papaspiliopoulos, 2018) that suggests to perform a Taylor expansion of  $\mathbf{p}_i(y_i|\psi_i)$  around the current state of the chain, leaving  $\mathbf{p}_i(\psi_i)$  unchanged.

## 5. Maximum Likelihood Estimation

### 5.1. The SAEM Algorithm

The ML estimator defined by (12) is computed using the Stochastic Approximation of the EM algorithm (SAEM) (Delyon et al., 1999). The SAEM algorithm is described as follows:

---

**Algorithm 3** The SAEM algorithm

---

**Initialization:**  $\theta_0$ , an initial parameter estimate and  $M$ , the number of MCMC iterations.

**Iteration k:** given the current model parameter estimate  $\theta_{k-1}$ :

1. **Simulation step:** for  $i \in \llbracket 1, N \rrbracket$ , draw vectors of individual parameters  $(\psi_1^{(k)}, \dots, \psi_N^{(k)})$  after  $M$  transitions of Markov kernels  $(\Pi_1^{(k)}(\psi_1^{(k-1)}, \cdot), \dots, (\Pi_N^{(k)}(\psi_N^{(k-1)}, \cdot)))$  which admit as unique limiting distributions the conditional distributions  $(\mathbf{p}_1(\psi_1|y_1; \theta_{k-1}), \dots, \mathbf{p}_N(\psi_N|y_N; \theta_{k-1}))$ ,
2. **Stochastic approximation step:** update the approximation of the conditional expectation  $\mathbb{E}[\log \mathbf{p}(y, \psi; \theta) | y, \theta_{k-1}]$ :

$$Q_k(\theta) = Q_{k-1}(\theta) + \gamma_k \left( \sum_{i=1}^N \log \mathbf{p}_i(y_i, \psi_i^{(k)}; \theta) - Q_{k-1}(\theta) \right), \quad (31)$$

where  $\{\gamma_k\}_{k>0}$  is a sequence of decreasing stepsizes with  $\gamma_1 = 1$ .

3. **Maximisation step:** Update the model parameter estimate:

$$\theta_k = \arg \max_{\theta \in \mathbb{R}^d} Q_k(\theta). \quad (32)$$


---

The SAEM algorithm is implemented in most software tools for NLME models and its convergence is studied in (Delyon et al., 1999; Kuhn and Lavielle, 2004; Allasonniere and Kuhn, 2013). The practical performances of SAEM are closely linked to the settings of SAEM. In particular, the choice of the transition kernel  $\Pi$  plays a key role. The transition kernel  $\Pi$  is directly defined by the proposal(s) used for the MH algorithm.

### 5.2. The f-SAEM algorithm

We propose a fast version of the SAEM algorithm using our resulting independent proposal distribution called the f-SAEM. The simulation step of the f-SAEM is achieved

using the nlme-IMH algorithm (see algorithm 2) for all individuals  $i \in \llbracket 1, N \rrbracket$  and the next steps remain unchanged.

**Remarks:**

1. This strongly relates to MAP algorithms (McLachlan and Krishnan, 2007) used for maximum likelihood estimation. Though, our f-SAEM algorithm consists in adding a rejection step on all the individual MAP estimates to bypass local trap issues.

In practice, the number of transitions  $M$  is small since the convergence of the SAEM does not require the convergence of the MCMC at each iteration (Kuhn and Lavielle, 2004). In the sequel, we carry out numerous numerical experiments to compare our nlme-IMH algorithm to state-of-the-art samplers and assess its relevance in a MLE algorithm such as the SAEM.

## 6. Numerical Examples

### 6.1. A pharmacokinetic example

#### 6.1.1. Data and model

32 healthy volunteers received a 1.5 mg/kg single oral dose of warfarin, an anticoagulant normally used in the prevention of thrombosis (O'Reilly and Aggeler, 1968). Figure 1 shows the warfarin plasmatic concentration measured at different times for these patients (the single dose was given at time 0 for all the patients).

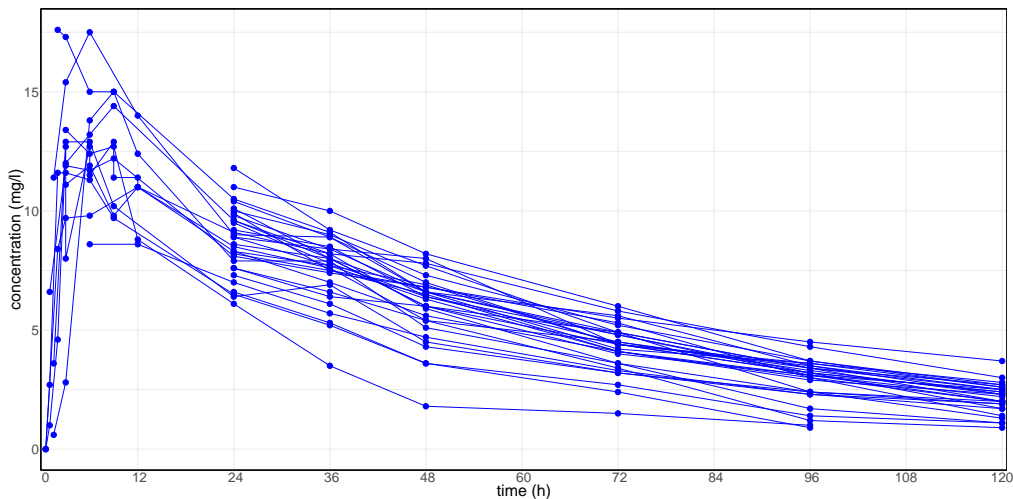


Figure 1: Warfarin concentration (mg/l) over time (h) for 32 subjects

We consider a one-compartment pharmacokinetics (PK) model for oral administration, assuming first-order absorption and linear elimination processes:

$$f(t, ka, V, k) = \frac{D ka}{V(ka - k)}(e^{-ka t} - e^{-k t}), \quad (33)$$

where  $ka$  is the absorption rate constant,  $V$  the volume of distribution,  $k$  the elimination rate constant, and  $D$  the dose of drug administered.

Here,  $ka$ ,  $V$  and  $k$  are PK parameters that can change from one individual to another. Let  $\psi_i = (ka_i, V_i, k_i)$  be the vector of individual PK parameters for individual  $i$ . The model for the  $j$ -th measured concentration, noted  $y_{ij}$ , for individual  $i$  writes:

$$y_{ij} = f(t_{ij}, \psi_i) + \varepsilon_{ij}. \quad (34)$$

We assume in this example that the residual errors are independent and normally distributed with mean 0 and variance  $\sigma^2$ . Lognormal distributions are used for the three PK parameters:

$$\log(ka_i) \sim \mathcal{N}(\log(ka_{\text{pop}}), \omega_{ka}^2), \quad (35)$$

$$\log(V_i) \sim \mathcal{N}(\log(V_{\text{pop}}), \omega_V^2), \quad (36)$$

$$\log(k_i) \sim \mathcal{N}(\log(k_{\text{pop}}), \omega_k^2). \quad (37)$$

This is a specific instance of the nonlinear mixed effects model for continuous data described in Section 2.2. We thus use the multivariate Gaussian proposal whose mean and covariance are defined by (23) and (21). In such case the gradient can be explicitly computed. Nevertheless, for the method to be easily extended to any structural model, the gradient is calculated using Automatic Differentiation (Griewank and Walther, 2008) implemented in the R package “Madness” (Pav, 2016).

### 6.1.2. MCMC Convergence Diagnostic

We study in this section the behaviour of the MH algorithm used to sample individual parameters from the conditional distribution  $\mathbf{p}_i(\psi_i|y_i; \theta)$ . We consider only one of the 32 individuals for this study and fix  $\theta$  to some arbitrary value, close to the ML estimate obtained with the SAEM algorithm, implemented in the saemix R package (Comets et al., 2017):  $ka_{\text{pop}} = 1$ ,  $V_{\text{pop}} = 8$ ,  $k_{\text{pop}} = 0.01$ ,  $\omega_{ka} = 0.5$ ,  $\omega_V = 0.2$ ,  $\omega_k = 0.3$  and  $\sigma^2 = 0.5$ .

We run the classical version of MH implemented in the saemix package and for which different transition kernels are used successively at each iteration: independent proposals from the marginal distribution  $\mathbf{p}_i(\psi_i)$ , component-wise random walk and block-wise random walk. We compare it to our proposed algorithm 2.

We run 20000 iterations of these two algorithms and evaluate their convergence by looking at the convergence of the empirical quantiles of order 0.1, 0.5 and 0.9 for the three components of  $\psi_i$ . Here,  $\hat{q}_\alpha^{(k)}(\psi_{i,\ell})$  is the empirical quantile of order  $\alpha$  of  $(\psi_{i,\ell}^{(1)}, \psi_{i,\ell}^{(2)}, \dots, \psi_{i,\ell}^{(k)})$  and  $\ell \in \llbracket 1, 3 \rrbracket$  denotes the component of the individual parameter.

We see Figure 2 that, for all  $\alpha$  and all  $\ell$ , the sequences of empirical quantiles  $\hat{q}_\alpha^k(\psi_{i,\ell})$  obtained with the two algorithms converge to the same value, which is supposed to be the theoretical quantile of the conditional distribution.

The interest of the new proposal is clearly visible here since all the empirical quantiles obtained with the nlme-IMH converge faster than with the reference algorithm.

Finally, it is interesting to note that the empirical medians converge very rapidly. This is interesting in the population approach framework because it is mainly the median values of each conditional distribution that are used to infer the population distribution.

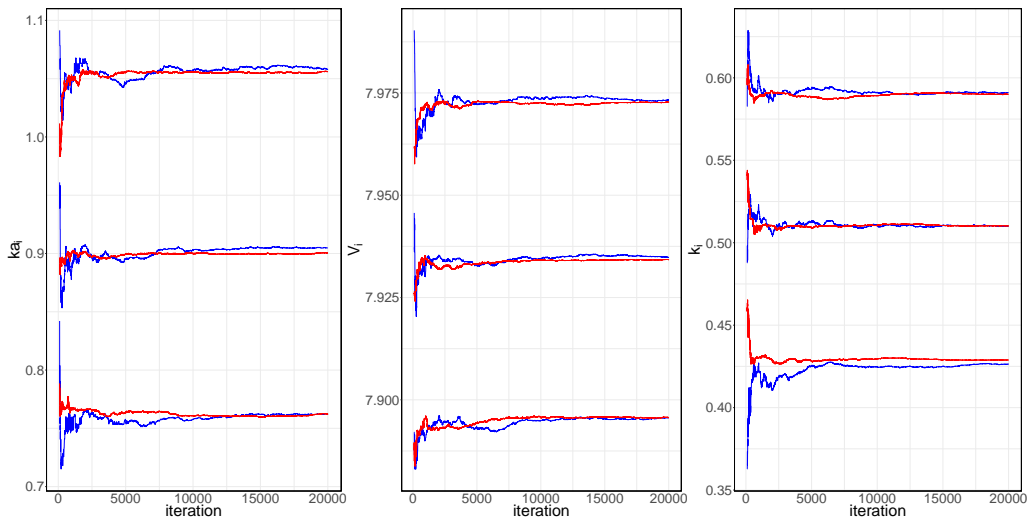


Figure 2: Modelling of the warfarin PK data: convergence of the empirical quantiles of order 0.1, 0.5 and 0.9 of  $\mathbf{p}_i(\psi_i|y_i; \theta)$  for a single individual. Comparison between the reference MH algorithm (blue) and the nlme-IMH (red).

*Comparison with state-of-the-art methods:* We then compare our new approach to the three following samplers: an independent sampler that uses variational approximation as proposal distribution (de Freitas et al., 2001), the MALA (Roberts and Tweedie, 1996) and the No-U-Turn Sampler (Hoffman and Gelman, 2014).

The same design and settings (dataset, model parameter estimate, individual) as in section 6.1.2 are used throughout the following experiments.

### 6.1.2.a. Variational MCMC algorithm

The Variational MCMC algorithm (de Freitas et al., 2001) is a MCMC algorithm with independent proposal. The proposal distribution is a multivariate Gaussian distribution whose parameters are obtained by a variational approach that consists in minimising the Kullback Leibler divergence between a multivariate Gaussian distribution  $q(\psi_i, \delta)$ , and the target distribution for a given model parameter estimate  $\theta$  noted  $\mathbf{p}_i(\psi_i|y_i, \theta)$ . This



problem boils down to maximizing the so-called Evidence Lower Bound  $\text{ELBO}(\theta)$  defined as:

$$\text{ELBO}(\delta) \triangleq \int q(\psi_i, \delta) (\log \mathbf{p}_i(y_i, \psi_i, \theta) - \log q(\psi_i, \delta)) d\psi_i \quad (38)$$

We use the Automatic Differentiation Variational Inference (ADVI) (Kucukelbir et al., 2015) implemented in RStan (R Package (Stan Development Team, 2018)) to obtain the vector of parameters noted  $\delta_{VI}$  defined as:

$$\delta_{VI} \triangleq \arg \max_{\delta \in \mathbb{R}^p \times \mathbb{R}^{p \times p}} \text{ELBO}(\delta) .$$

The algorithm stops when the variation of the median of the objective function falls below the 1% threshold. The means of our nlme-IMH and the Variational MCMC proposals compare with the posterior mean (calculated using the NUTS (Hoffman and Gelman, 2014)) as follows:

Table 1: Means of the proposals.

	$ka_i$	$V_i$	$k_i$
Variational MCMC	0.90	7.93	0.48
<b>nlme-IMH</b>	0.88	7.93	0.52
NUTS	0.91	7.93	0.51

We see Figure 3 that the independent sampler using a variational approximation as a proposal performs poorly in this example.

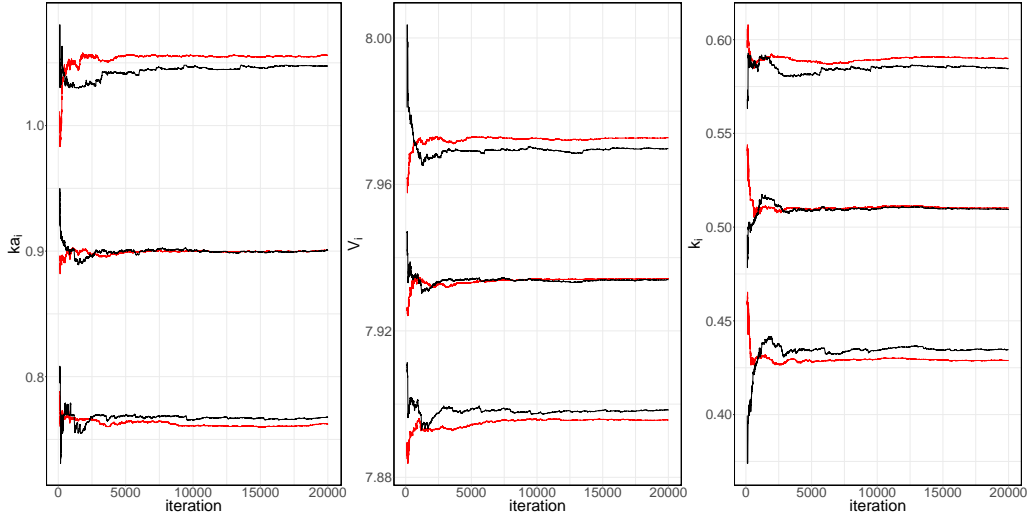


Figure 3: Modelling of the warfarin PK data: convergence of the empirical quantiles of order 0.1, 0.5 and 0.9 of  $\mathbf{p}_i(\psi_i|y_i; \theta)$  for a single individual. Comparison between the nlme-IMH (red) and the Variational MCMC (black).

We observe that the mean of the variational approximation is slightly shifted from the estimated posterior mode (see table 1 for comparison) whereas a considerable difference lies in the numerical value of the covariance matrix obtained with ADVI. The empirical standard deviation of the Variational MCMC proposal is much smaller than our new proposal defined by (21) (see table 2), which slows down the MCMC convergence.

Table 2: Standard deviations of the proposals.

	$ka_i$	$V_i$	$k_i$
Variational MCMC	0.14	0.03	0.07
<b>nlme-IMH</b>	0.18	0.04	0.09
NUTS	0.18	0.05	0.09

Table 3: Pairwise correlations of the proposals.

	$ka_i, V_i$	$ka_i, k_i$	$V_i, k_i$
Variational MCMC	0.48	-0.28	-0.61
<b>nlme-IMH</b>	0.56	-0.39	-0.68
NUTS	0.55	-0.39	-0.68

Figure 4 shows the proposals marginals and the marginal posterior distribution for the individual parameters  $k_i$  and  $V_i$ . Biplot of the samples drawn from the two multivariate Gaussian proposals (our independent proposal and the variational MCMC proposal) as well as samples drawn from the posterior distribution (using the NUTS) are also presented in this figure. We conclude that both marginal and bivariate posterior distributions are better approximated by our independent proposal than the one resulting from a KL divergence optimization.

Besides similar marginal variances, both our independent proposal and the true posterior share a strong anisotropic nature, confirmed by the similar correlation values of table 3. Same characteristics are observed for the other parameters.

### 6.1.2.b. Metropolis Adjusted Langevin Algorithm (MALA)

We now compare our method to the MALA, which proposal is defined by (15). The gradient of the log posterior distribution  $\nabla_{\psi_i} \log p_i(\psi_i^{(k)} | y_i)$  is also calculated by Automatic Differentiation. In this numerical example, the MALA has been initialized at the MAP and the stepsize ( $\gamma = 10^{-2}$ ) is tuned such that the acceptance rate of 0.57 is reached (Roberts and Rosenthal, 1997).

Figure 5 highlights good convergence of a well-tuned MALA. Quantiles stabilisation is quite similar to our method for all orders and components.

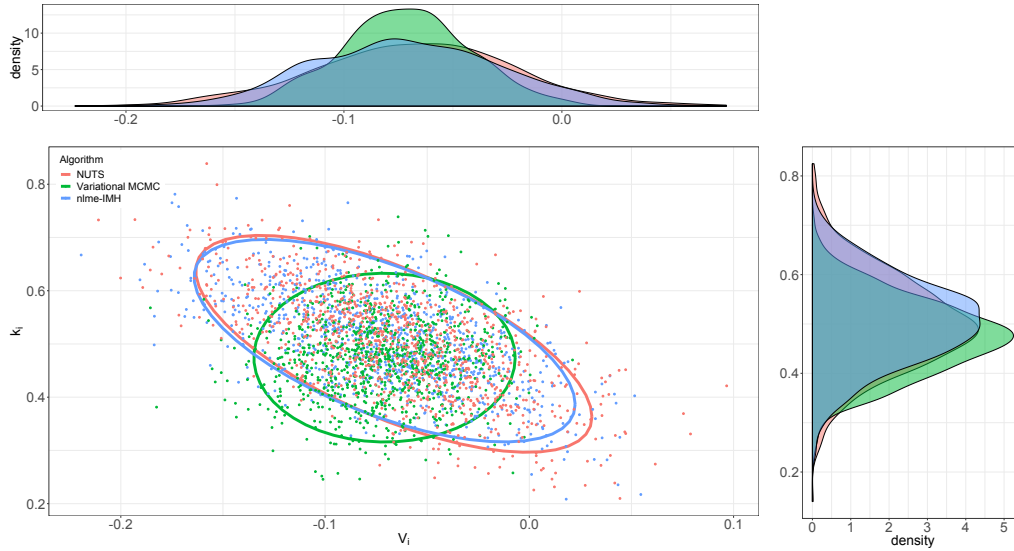


Figure 4: Modelling of the warfarin PK data: Comparison between the proposals of the nlme-IMH (blue), the Variational MCMC (green) and the empirical target distribution sampled using the NUTS (red). Marginals and biplots of the conditional distributions  $k_i|y_i$  and  $V_i|y_i$  for a single individual. Ellipses containing 90% of the data points are represented on the main plot.

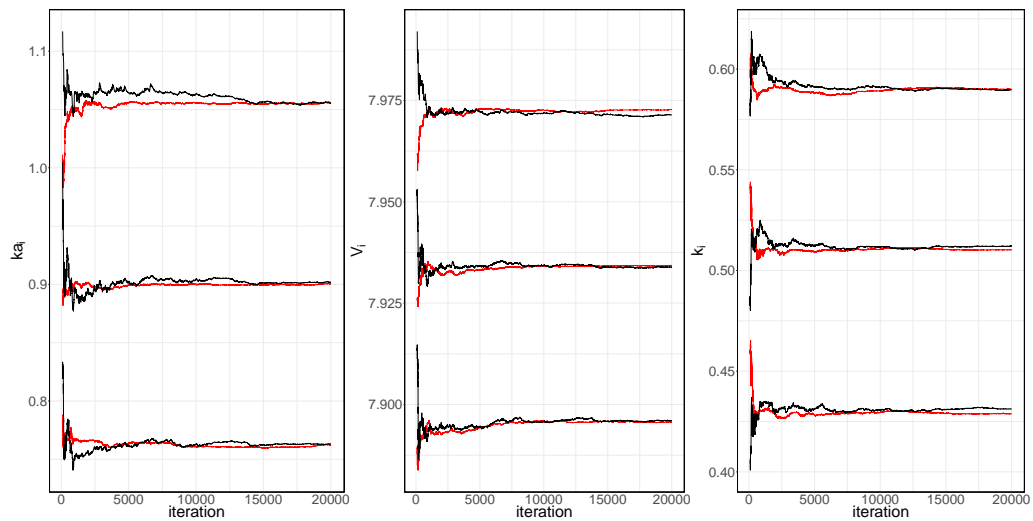


Figure 5: Modelling of the warfarin PK data: convergence of the empirical quantiles of order 0.1, 0.5 and 0.9 of  $p_i(\psi_i|y_i; \theta)$  for a single individual. Comparison between the nlme-IMH (red) and the MALA (black).

### 6.1.2.c. No-U-Turn Sampler (NUTS)

We compare the implementation of NUTS (Hoffman and Gelman, 2014; Carpenter et al., 2017) in the RStan package to our method in Figure 6. We observe that the empirical quantiles obtained with the NUTS steadily converge to the target values.

Even though the behaviour of our method seems to be similar in the long run, in the first 1 000 the nlme-IMH algorithm stabilizes a bit more slowly than the NUTS around a neighbourhood of the target values.

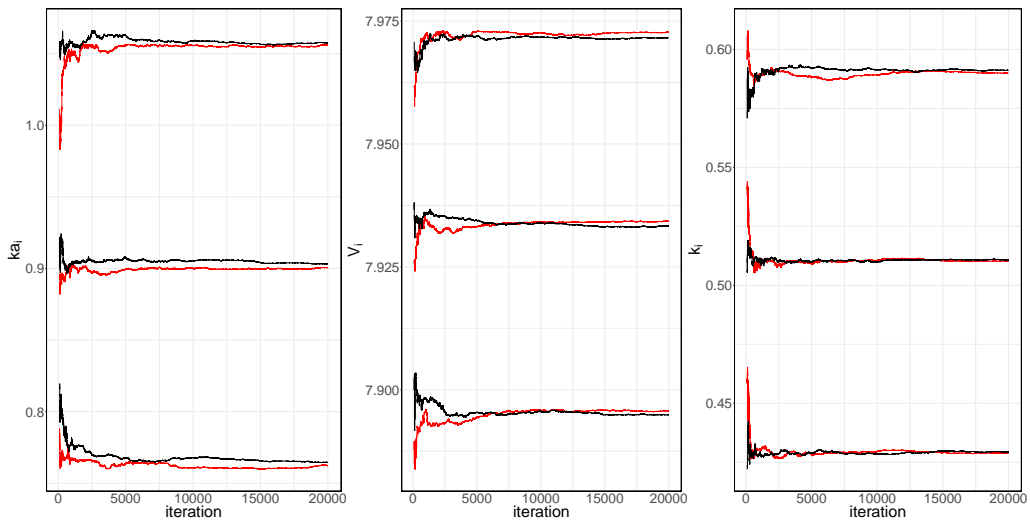


Figure 6: Modelling of the warfarin PK data: convergence of the empirical quantiles of order 0.1, 0.5 and 0.9 of  $p_i(\psi_i|y_i; \theta)$  for a single individual. Comparison between the nlme-IMH (red) and the NUTS (black).

In practice, those three methods imply tuning phases that are computationally involved, warming up the chain and a careful initialisation whereas our independent sampler is automatic and fast to implement. Investigating the asymptotic convergence behavior of those methods highlights the competitive properties of our IMH algorithm to sample from the target distribution.

Since our goal is to embed those samplers into a MLE algorithm such as the SAEM, we shall now study how they behave in the very first iterations of the MCMC procedure. Recall that the SAEM requires only few iterations of MCMC sampling under the current model parameter estimate. We present this non asymptotic study in the following section.

#### 6.1.3. Comparison of the chains for the first 100 iterations

We produce 100 independent runs of the RWM, the nlme-IMH, the MALA and the NUTS for 500 iterations. The boxplots of the samples drawn at a given iteration threshold (three different thresholds are used) are presented Figure 7 against the ground truth

for the parameter **ka**. The ground truth has been calculated by running the NUTS for 100 000 iterations.

For the three numbers of iteration (5,20,500) considered in Figure 7, the median of the nlme-IMH and NUTS samples are closer to the ground truth. Figure 7 also highlights that all those methods succeed in sampling from the whole distribution after 500 iterations.

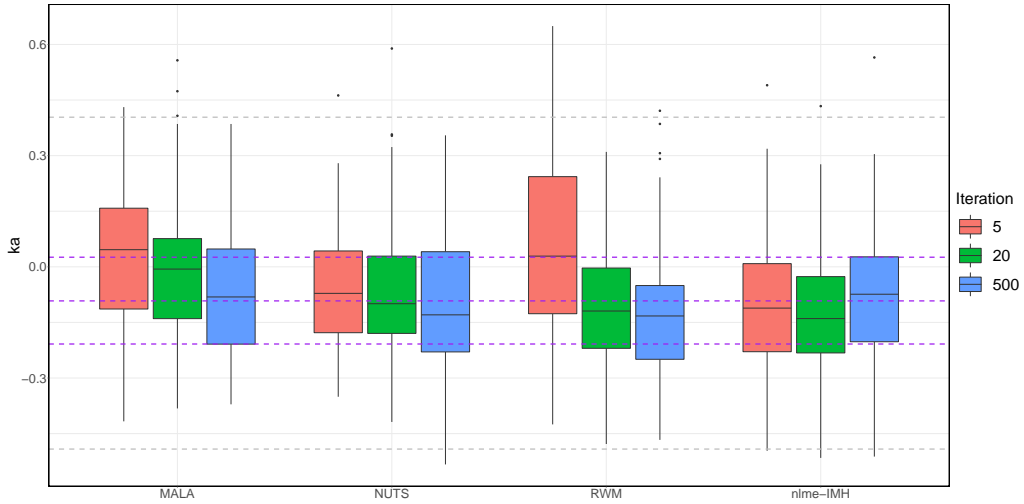


Figure 7: Modelling of the warfarin PK data: Boxplots for the RWM, the nlme-IMH, the MALA and the NUTS algorithm, averaged over 100 independent runs. The groundtruth median, 0.25 and 0.75 percentiles are plotted as a dashed purple line and its maximum and minimum as a dashed grey line.

We now use the RWM, the nlme-IMH and the MALA in the SAEM algorithm and observe the convergence of the resulting sequences of parameters.

#### 6.1.4. Maximum likelihood estimation

We use the SAEM algorithm to estimate the population PK parameters  $ka_{\text{pop}}$ ,  $V_{\text{pop}}$  and  $k_{\text{pop}}$ , the standard deviations of the random effects  $\omega_{ka}$ ,  $\omega_V$  and  $\omega_k$  and the residual variance  $\sigma^2$ .

The stepsize  $\gamma_k$  is set to 1 during the first 100 iterations and then decreases as  $1/k^a$  where  $a = 0.7$  during the next 100 iterations.

Here we compare the standard SAEM algorithm, as implemented in the saemix R package, with the f-SAEM algorithm and the SAEM using the MALA sampler. In this example, the nlme-IMH and the MALA are only used during the first 20 iterations of the SAEM. The standard MH algorithm is then used.

Figure 8 shows the estimates of  $V_{\text{pop}}$  and  $\omega_V$  computed at each iteration of these three variants of SAEM and starting from three different initial values. First of all,

we notice that, whatever the initialisation and the sampling algorithm used, all the runs converge towards the maximum likelihood estimate. It is then very clear that the f-SAEM converges faster than the standard algorithm. The SAEM using the MALA algorithm for sampling from the individual conditional distribution presents a similar convergence behavior as the reference.

We can conclude, for this example, that sampling around the MAP of each individual conditional distribution is the key to a fast convergence of the SAEM during the first iterations.

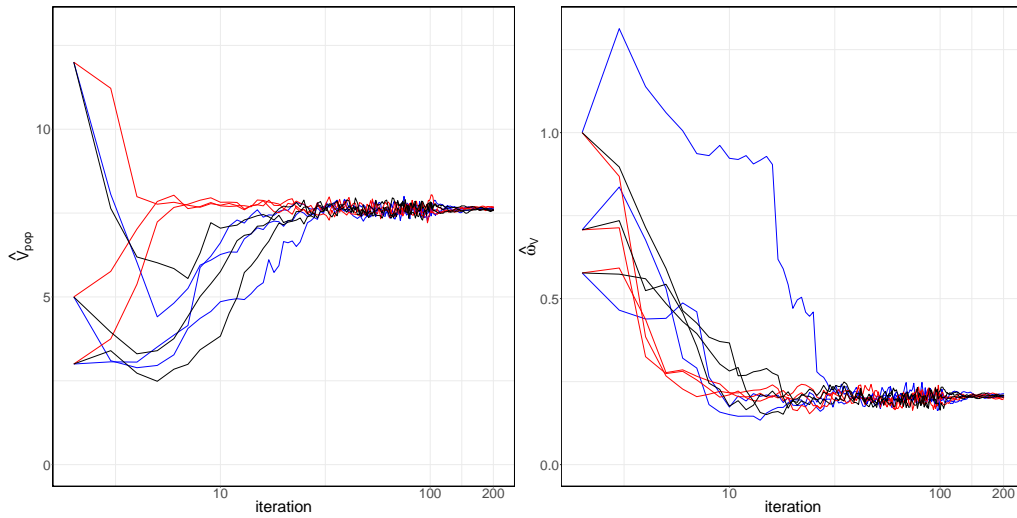


Figure 8: Estimation of the population PK parameters for the warfarin data: convergence of the sequences of estimates  $(\hat{V}_{\text{pop},k}, 1 \leq k \leq 200)$  and  $(\hat{\omega}_{V,k}, 1 \leq k \leq 200)$  obtained with SAEM and three different initial values using the reference MH algorithm (blue), the f-SAEM (red) and the SAEM using the MALA sampler (black).

### 6.1.5. Monte Carlo study

We conduct a Monte Carlo study to confirm the properties of the f-SAEM algorithm for computing the ML estimates.

$M = 50$  datasets have been simulated using the PK model previously used for fitting the warfarin PK data with the following parameter values:  $ka_{\text{pop}} = 1$ ,  $V_{\text{pop}} = 8$ ,  $k_{\text{pop}} = 0.1$ ,  $\omega_{ka} = 0.5$ ,  $\omega_V = 0.2$ ,  $\omega_k = 0.3$  and  $\sigma^2 = 0.5$ . The same original design with  $N = 32$  patients and a total number of 251 PK measurements were used for all the simulated datasets.

Since all the simulated data are different, the value of the ML estimator varies from one simulation to another. If we run  $K$  iterations of SAEM, the last element of the sequence  $(\theta_k^{(m)}, 1 \leq k \leq K)$  is the estimate obtained from the  $m$ -th simulated dataset. To investigate how fast  $(\theta_k^{(m)}, 1 \leq k \leq K)$  converges to  $\theta_K^{(m)}$  we study how fast  $(\theta_k^{(m)} -$

$\theta_K^{(m)}, 1 \leq k \leq K$ ) goes to 0.

For a given sequence of estimates, we can then define, at each iteration  $k$  and for each component  $\ell$  of the parameter, the mean square distance over the replicates

$$E_k(\ell) = \frac{1}{M} \sum_{m=1}^M \left( \theta_k^{(m)}(\ell) - \theta_K^{(m)}(\ell) \right)^2. \quad (39)$$

Figure 9 shows using the new proposal leads to a much faster convergence towards the maximum likelihood estimate. Less than 10 iterations are required to converge with the f-SAEM on this example, instead of 50 with the original version. It should also be noted that the distance decreases monotonically. The sequence of estimates approaches the target at each iteration, compared to the standard algorithm which makes twists and turns before converging.

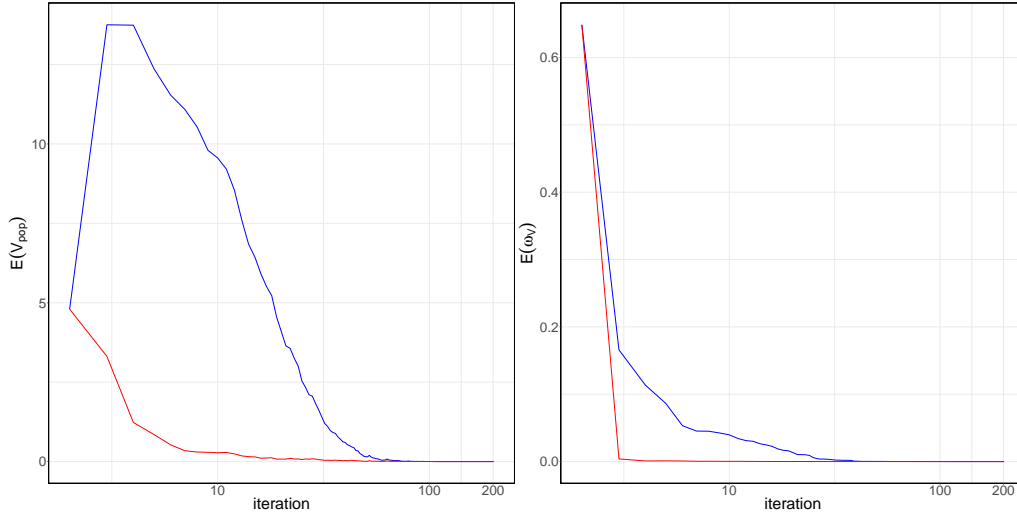


Figure 9: Convergence of the sequences of mean square distances  $(E_k(V_{\text{pop}}), 1 \leq k \leq 200)$  and  $(E_k(\omega_V), 1 \leq k \leq 200)$  for  $V_{\text{pop}}$  and  $\omega_V$  obtained with SAEM on  $M = 50$  synthetic datasets using the reference MH algorithm (blue) and the f-SAEM (red).

## 6.2. Time-to-event Data Model

### 6.2.1. The model

In this section, we consider a Weibull model for time-to-event data (Lavielle, 2014; Zhang, 2016). For individual  $i$ , the hazard function of this model is:

$$h(t, \psi_i) = \frac{\beta_i}{\lambda_i} \left( \frac{t}{\lambda_i} \right)^{\beta_i - 1}. \quad (40)$$

Here, the vector of individual parameters is  $\psi_i = (\lambda_i, \beta_i)$ . These two parameters are assumed to be independent and lognormally distributed:

$$\log(\lambda_i) \sim \mathcal{N}(\log(\lambda_{\text{pop}}), \omega_\lambda^2), \quad (41)$$

$$\log(\beta_i) \sim \mathcal{N}(\log(\beta_{\text{pop}}), \omega_\beta^2). \quad (42)$$

Then, the vector of population parameters is  $\theta = (\lambda_{\text{pop}}, \beta_{\text{pop}}, \omega_\lambda, \omega_\beta)$ .

Repeated events were generated, for  $N = 100$  individuals, using the Weibull model (40) with  $\lambda_{\text{pop}} = 10$ ,  $\omega_\lambda = 0.3$ ,  $\beta_{\text{pop}} = 3$  and  $\omega_\beta = 0.3$  and assuming a right censoring time  $\tau_c = 20$ .

### 6.2.2. MCMC Convergence Diagnostic

Similarly to the previous section, we start by looking at the behaviour of the MCMC procedure used for sampling from the conditional distribution  $\mathbf{p}_i(\psi_i|y_i; \theta)$  for a given individual  $i$  and assuming that  $\theta$  is known. We use the generating model parameter in these experiments ( $\theta = (\lambda_{\text{pop}} = 10, \beta_{\text{pop}} = 3, \omega_\lambda = 0.3, \omega_\beta = 0.3)$ ).

We ran 10 000 iterations of the reference MH algorithm the nlme-IMH to estimate quantiles of order 0.1, 0.5 and 0.9 of the conditional distributions of  $\lambda_i$  and  $\beta_i$ .

We see Figure 10 that the sequences of empirical quantiles obtained with the two procedures converge to the same value but the new algorithm converges faster than the standard MH algorithm.

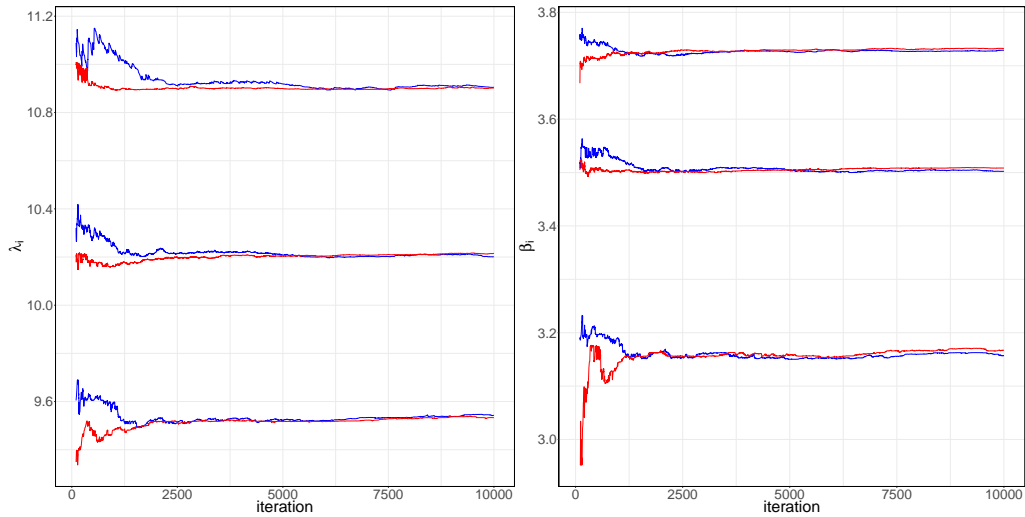


Figure 10: Time-to-event data modelling: convergence of the empirical quantiles of order 0.1, 0.5 and 0.9 of  $\mathbf{p}_i(\psi_i|y_i; \theta)$  for a single individual. The reference MH algorithm is in blue and the nlme-IMH is in red.



Comparisons with state-of-the-art methods were conducted as in the previous section. These comparisons led us to the same remarks as those made for the previous continuous data model both on the asymptotic and non asymptotic regimes.

### 6.2.3. Maximum likelihood estimation of the population parameters

We run the standard SAEM algorithm implemented in the saemix package (extension of this package for noncontinuous data models is available on GitHub: <https://github.com/belhal/saemix>) and the f-SAEM on the generated dataset.

Figure 11 shows the estimates of  $\lambda_{\text{pop}}$  and  $\omega_\lambda$  computed at each iteration of the two versions of the SAEM and starting from three different initial values. The same behaviour is observed as in the continuous case: regardless the initial values and the algorithm, all the runs converge to the same solution but convergence is much faster with the proposed method. The same comment applies for the two other parameters  $\beta_{\text{pop}}$  and  $\omega_\beta$ .

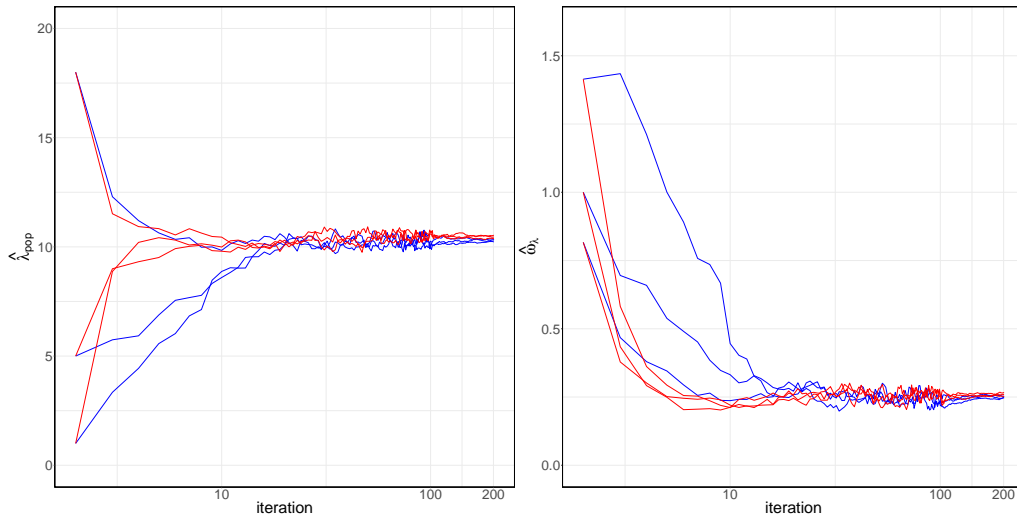


Figure 11: Population parameter estimation in time-to-event-data models: convergence of the sequences of estimates  $(\hat{\lambda}_{\text{pop},k}, 1 \leq k \leq 200)$  and  $(\hat{\omega}_{\lambda,k}, 1 \leq k \leq 200)$  obtained with SAEM and three different initial values using the reference MH algorithm (blue) and the f-SAEM (red).

### 6.2.4. Monte Carlo study

We now conduct a Monte Carlo study in order to confirm the good properties of the new version of the SAEM algorithm for estimating the population parameters of a time-to-event data model.

$M = 50$  synthetic datasets are generated using the same design as above. Figure 12 shows the convergence of the mean square distances defined in (39) for  $\lambda_{\text{pop}}$  and  $\omega_\lambda$ . All these distances converge monotonically to 0 which means that both algorithms properly

converge to the maximum likelihood estimate, but very few iterations are required with the new version to converge while about thirty iterations are needed with the standard SAEM.

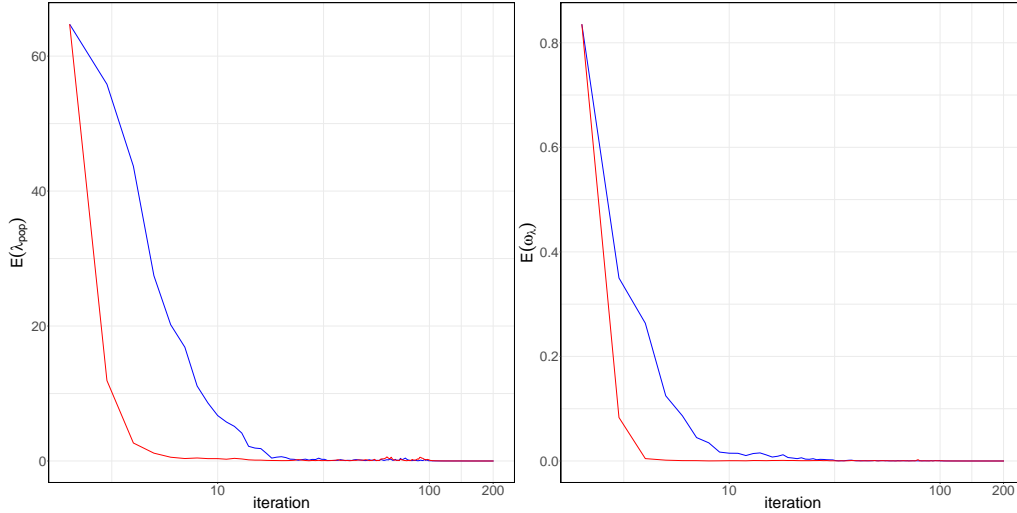


Figure 12: Convergence of the sequences of mean square distances  $(E_k(\lambda_{pop}), 1 \leq k \leq 200)$  and  $(E_k(\omega_\lambda), 1 \leq k \leq 200)$  for  $\lambda_{pop}$  and  $\omega_\lambda$  obtained with SAEM from  $M = 50$  synthetic datasets using the reference MH algorithm (blue) and the f-SAEM (red).

## 7. Conclusion and discussion

We present in this article an independent Metropolis-Hastings procedure for sampling random effects from their conditional distributions in nonlinear mixed effects models.

The idea of the method is to approximate each individual conditional distribution by a multivariate normal distribution. A Laplace approximation makes it possible to consider any type of data, but we have shown that, in the case of continuous data, this approximation is equivalent to linearizing the structural model around the conditional mode of the random effects.

The numerical experiments demonstrate that the proposed nlme-IMH sampler converges faster to the target distribution than a standard random walk Metropolis. This practical behaviour is partly explained by the fact that the conditional mode of the random effects in the linearized model coincides with the conditional mode of the random effects in the original model. The proposal distribution is therefore a normal distribution centered around this MAP. On the other hand, the dependency structure in the conditional distribution of the random effects is well approximated by the covariance structure of the Gaussian proposal.

Comparison studies between our approach and the independent sampler using a variational approximation proposal, the MALA and the NUTS have shown similar quantile convergence behaviour.

So far, we have mainly applied our method to standard problems encountered in pharmacometrics and for which the number of random effects remains small. It can nevertheless be interesting to see how this method behaves in higher dimension and compare it with methods adapted to such situations such as MALA or HMC. Lastly, we have shown that this new IMH algorithm can easily be embedded in the SAEM algorithm for maximum likelihood estimation of the population parameters. Our numerical studies have shown empirically that the new transition kernel is effective in the very first iterations. It is of great interest to determine automatically and in an adaptive way an optimal scheme of kernel transitions combining this new proposal with the block-wise random walk Metropolis.

## References

- Agresti, A., 1990. Categorical data analysis. A Wiley-Interscience publication, Wiley, New York.
- Allasonniere, S., Kuhn, E., 2013. Convergent Stochastic Expectation Maximization algorithm with efficient sampling in high dimension. Application to deformable template model estimation. arXiv preprint arXiv:1207.5938 .
- Andersen, P.K., 2006. Survival Analysis. Wiley Reference Series in Biostatistics .
- Andrieu, C., Moulines, É., et al., 2006. On the ergodicity properties of some adaptive mcmc algorithms. *The Annals of Applied Probability* 16, 1462–1505.
- Andrieu, C., Thoms, J., 2008. A tutorial on adaptive mcmc. *Statistics and computing* 18, 343–373.
- Atchadé, Y.F., Rosenthal, J.S., 2005. On adaptive Markov chain Monte Carlo algorithms. *Bernoulli* 11, 815–828. doi:[10.3150/bj/1130077595](https://doi.org/10.3150/bj/1130077595).
- Beal, S., Sheiner, L., 1980. The NONMEM system. *The American Statistician* 34, 118–119.
- Betancourt, M., 2017. A conceptual introduction to Hamiltonian Monte Carlo. arXiv preprint arXiv:1701.02434 .
- Brooks, S., Gelman, A., Jones, G., Meng, X.L., 2011. Handbook of markov chain monte carlo. CRC press.
- Brosse, N., Durmus, A., Moulines, É., Sabanis, S., 2017. The tamed unadjusted langevin algorithm. arXiv preprint arXiv:1710.05559 .
- Carpenter, B., Gelman, A., Hoffman, M., Lee, D., Ben, G., Betancourt, M., Brubaker, M., Guo, J., Li, P., Riddell, A., 2017. Stan: A probabilistic programming language. *Journal of Statistical Software* 76.
- Chan, P.L.S., Jacqmin, P., Lavielle, M., McFadyen, L., Weatherley, B., 2011. The use of the SAEM algorithm in MONOLIX software for estimation of population pharmacokinetic-pharmacodynamic-viral dynamics parameters of maraviroc in asymptomatic HIV subjects. *Journal of Pharmacokinetics and Pharmacodynamics* 38, 41–61.
- Comets, E., Lavenu, A., Lavielle, M., 2017. Parameter estimation in nonlinear mixed effect models using saemix, an r implementation of the saem algorithm. *Journal of Statistical Software* 80, 1–42.
- Delyon, B., Lavielle, M., Moulines, E., 1999. Convergence of a stochastic approximation version of the EM algorithm. *Ann. Statist.* 27, 94–128. doi:[10.1214/aos/1018031103](https://doi.org/10.1214/aos/1018031103).
- Durmus, A., Roberts, G.O., Vilmart, G., Zygalakis, K.C., 2017. Fast langevin based algorithm for mcmc in high dimensions. *Ann. Appl. Probab.* 27, 2195–2237. doi:[10.1214/16-AAP1257](https://doi.org/10.1214/16-AAP1257).
- de Freitas, N., Høj-Sørensen, P., Jordan, M.I., Russell, S., 2001. Variational mcmc. *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence* , 120–127.
- Girolami, M., Calderhead, B., 2011. Riemann manifold langevin and hamiltonian monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73, 123–214.
- Griewank, A., Walther, A., 2008. Evaluating derivatives: principles and techniques of algorithmic differentiation. volume 105. Siam.
- Haario, H., Saksman, E., Tamminen, J., et al., 2001. An adaptive metropolis algorithm. *Bernoulli* 7, 223–242.
- Hoffman, M.D., Gelman, A., 2014. The No-U-turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research* 15, 1593–1623.
- Kucukelbir, A., Ranganath, R., Gelman, A., Blei, D., 2015. Automatic variational inference in stan, in: Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems* 28. Curran Associates, Inc., pp. 568–576.
- Kuhn, E., Lavielle, M., 2004. Coupling a stochastic approximation version of EM with an MCMC procedure. *ESAIM: Probability and Statistics* 8, 115–131.
- Lavielle, M., 2014. Mixed effects models for the population approach: models, tasks, methods and tools. CRC press.
- Lavielle, M., Ribba, B., 2016. Enhanced method for diagnosing pharmacometric models: random sampling from conditional distributions. *Pharmaceutical research* 33, 2979–2988.
- Louis, T.A., 1982. Finding the observed information matrix when using the em algorithm. *Journal of the Royal Statistical Society, Series B: Methodological* 44, 226–233.
- Mbogning, C., Bleakley, K., Lavielle, M., 2015. Joint modeling of longitudinal and repeated time-to-event data using nonlinear mixed-effects models and the SAEM algorithm. *Journal of Statistical Computation and Simulation* 85, 1512–1528. doi:[10.1080/00949655.2013.878938](https://doi.org/10.1080/00949655.2013.878938).
- McLachlan, G., Krishnan, T., 2007. The EM algorithm and extensions. volume 382. John Wiley & Sons.
- Mengersen, K.L., Tweedie, R.L., 1996. Rates of convergence of the hastings and metropolis algorithms. *Ann. Statist.* 24, 101–121. doi:[10.1214/aos/1033066201](https://doi.org/10.1214/aos/1033066201).

- Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., Teller, E., 1953. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics* 21, 1087–1092. doi:[10.1063/1.1699114](https://doi.org/10.1063/1.1699114).
- Migon, H., Gamerman, D., Louzada, F., 2014. *Statistical Inference: An Integrated Approach*, Second Edition. Chapman & Hall/CRC Texts in Statistical Science, Taylor & Francis.
- Neal, R.M., et al., 2011. Mcmc using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo* 2, 2.
- O'Reilly, R.A., Aggeler, P.M., 1968. Studies on coumarin anticoagulant drugs initiation of warfarin therapy without a loading dose. *Circulation* 38, 169–177.
- Pav, S.E., 2016. Madness: a package for multivariate automatic differentiation .
- Robert, C.P., Casella, G., 2010. *Metropolis–Hastings Algorithms*. Springer New York, New York, NY. pp. 167–197. doi:[10.1007/978-1-4419-1576-4\\_6](https://doi.org/10.1007/978-1-4419-1576-4_6).
- Roberts, G.O., Gelman, A., Gilks, W.R., 1997. Weak convergence and optimal scaling of random walk metropolis algorithms. *Ann. Appl. Probab.* 7, 110–120. doi:[10.1214/aop/1034625254](https://doi.org/10.1214/aop/1034625254).
- Roberts, G.O., Rosenthal, J.S., 1997. Optimal scaling of discrete approximations to langevin diffusions. *J. R. Statist. Soc. B* 60, 255–268.
- Roberts, G.O., Rosenthal, J.S., 2011. Quantitative non-geometric convergence bounds for independence samplers. *Methodology and Computing in Applied Probability* 13, 391–403.
- Roberts, G.O., Tweedie, R.L., 1996. Exponential convergence of langevin distributions and their discrete approximations. *Bernoulli* 2, 341–363.
- Rue, H., Martino, S., Chopin, N., 2009. Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of the royal statistical society: Series b (statistical methodology)* 71, 319–392.
- Savic, R.M., Mentré, F., Lavielle, M., 2011. Implementation and evaluation of the SAEM algorithm for longitudinal ordered categorical data with an illustration in pharmacokinetics-pharmacodynamics. *The AAPS Journal* 13, 44–53.
- Stan Development Team, 2018. RStan: the R interface to Stan. URL: <http://mc-stan.org/>. r package version 2.17.3.
- Stramer, O., Tweedie, R.L., 1999. Langevin-type models i: Diffusions with given stationary distributions and their discretizations\*. *Methodology And Computing In Applied Probability* 1, 283–306. doi:[10.1023/A:1010086427957](https://doi.org/10.1023/A:1010086427957).
- Titsias, M.K., Papaspiliopoulos, O., 2018. Auxiliary gradient-based sampling algorithms. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 0. doi:[10.1111/rssb.12269](https://doi.org/10.1111/rssb.12269).
- Verbeke, G., 1997. *Linear mixed models for longitudinal data*. Springer.
- Vihola, M., 2012. Robust adaptive metropolis algorithm with coerced acceptance rate. *Statistics and Computing* 22, 997–1008.
- Zhang, Z., 2016. Parametric regression model for survival data: Weibull regression model as an example. *Ann Transl Med.* 24.

# Appendices

## A. Calculus of the proposal in the noncontinuous case

Laplace approximation (see (Migon et al., 2014)) consists in approximating an integral of the form

$$I := \int e^{v(x)} dx, \quad (43)$$

where  $v$  is at least twice differentiable.

The following second order Taylor expansion of the function  $v$  around a point  $x_0$

$$v(x) \approx v(x_0) + \nabla v(x_0)(x - x_0) + \frac{1}{2}(x - x_0)\nabla^2 v(x_0)(x - x_0), \quad (44)$$

provides an approximation of the integral  $I$  (consider a multivariate Gaussian probability distribution function which integral sums to 1):

$$I \approx e^{v(x_0)} \sqrt{\frac{(2\pi)^p}{|-\nabla^2 v(x_0)|}} \exp\left\{-\frac{1}{2}\nabla v(x_0)'\nabla^2 v(x_0)^{-1}\nabla v(x_0)\right\}. \quad (45)$$

In our context, we can write the marginal pdf  $\mathbf{p}_i(y_i)$  that we aim to approximate as

$$\mathbf{p}_i(y_i) = \int \mathbf{p}_i(y_i, \psi_i) d\psi_i \quad (46)$$

$$= \int e^{\log \mathbf{p}_i(y_i, \psi_i)} d\psi_i. \quad (47)$$

Then, let

$$v(\psi_i) := \log \mathbf{p}_i(y_i, \psi_i) \quad (48)$$

$$= l(\psi_i) + \log \mathbf{p}_i(\psi_i), \quad (49)$$

and compute its Taylor expansion around the MAP  $\hat{\psi}_i$ . We have by definition that

$$\nabla \log \mathbf{p}_i(y_i, \hat{\psi}_i) = 0,$$

which leads to the following Laplace approximation of  $\log \mathbf{p}_i(y_i)$ :

$$-2 \log \mathbf{p}_i(y_i) \approx -p \log 2\pi - 2 \log \mathbf{p}_i(y_i, \hat{\psi}_i) + \log \left( \left| -\nabla^2 \log \mathbf{p}_i(y_i, \hat{\psi}_i) \right| \right).$$

We thus obtain the following approximation of the logarithm of the conditional pdf of  $\psi_i$  given  $y_i$  evaluated at  $\hat{\psi}_i$ :

$$\log \mathbf{p}_i(\hat{\psi}_i | y_i) \approx -\frac{p}{2} \log 2\pi - \frac{1}{2} \log \left( \left| -\nabla^2 \log \mathbf{p}_i(y_i, \hat{\psi}_i) \right| \right),$$

which is precisely the log-pdf of a multivariate Gaussian distribution with mean  $\hat{\psi}_i$  and variance-covariance  $-\nabla^2 \log \mathbf{p}_i(y_i, \hat{\psi}_i)^{-1}$  with:

$$\begin{aligned} \nabla^2 \log \mathbf{p}_i(y_i, \hat{\psi}_i) &= \nabla^2 \log \mathbf{p}_i(y_i | \hat{\psi}_i) + \nabla^2 \log \mathbf{p}_i(\hat{\psi}_i) & (50) \\ &= \nabla^2 l(\hat{\psi}_i) + \Omega^{-1} . & (51) \end{aligned}$$

## B. Linear continuous data models

Let  $y_i = (y_{i,1}, \dots, y_{i,n_i})'$  and  $\varepsilon_i = (\varepsilon_{i,1}, \dots, \varepsilon_{i,n_i})'$ . Assume a linear relationship between the observations  $y_i$  and the vector of individual parameters  $\psi_i$ :

$$y_i = A_i \psi_i + \varepsilon_i , \quad (52)$$

where  $A_i \in \mathbb{R}^{n_i \times p}$  is the design matrix for individual  $i$ ,  $\psi_i$  is normally distributed with mean  $m_i \in \mathbb{R}^p$  and covariance  $\Omega \in \mathbb{R}^{p \times p}$ . Then, the conditional distribution of  $\psi_i$  given  $y_i$  is a normal distribution with mean  $\mu_i$  and variance-covariance matrix  $\Gamma_i$  defined as:

$$\mu_i = \Gamma_i \left( \frac{A_i' y_i}{\sigma^2} + \Omega^{-1} m_i \right) \quad \text{where} \quad \Gamma_i = \left( \frac{A_i' A_i}{\sigma^2} + \Omega^{-1} \right)^{-1} \quad (53)$$

Here,  $\mu_i$  is the mode of the conditional distribution of  $\psi_i$ , known as the Maximum A Posteriori (MAP) estimate, or the Empirical Bayes Estimate (EBE) of  $\psi_i$ .