

## ARTICLE

# SAMBA: A novel method for fast automatic model building in nonlinear mixed-effects models

Mélanie Prague<sup>1,2</sup> | Marc Lavielle<sup>3</sup>

<sup>1</sup>Inria Bordeaux Sud-Ouest, Inserm, Bordeaux Population Health Research Center, SISTM Team, UMR 1219, University of Bordeaux, Bordeaux, France

<sup>2</sup>Vaccine Research Institute, Créteil, France

<sup>3</sup>Inria & CMAP, Ecole Polytechnique, CNRS, Institut Polytechnique de Paris, Paris, France

**Correspondence**

Mélanie Prague, ISPED - Université de Bordeaux - Bureau 23, 146 Rue Léo Saignat, 33070 Bordeaux Cedex, France. Email: Melanie.Prague@inria.fr

**Funding information**

This study has received funding from the Nipah virus project financed by the French Ministry of Higher Education, Research, and Innovation

**Abstract**

The success of correctly identifying all the components of a nonlinear mixed-effects model is far from straightforward: it is a question of finding the best structural model, determining the type of relationship between covariates and individual parameters, detecting possible correlations between random effects, or also modeling residual errors. We present the Stochastic Approximation for Model Building Algorithm (SAMBA) procedure and show how this algorithm can be used to speed up this process of model building by identifying at each step how best to improve some of the model components. The principle of this algorithm basically consists in “learning something” about the “best model,” even when a “poor model” is used to fit the data. A comparison study of the SAMBA procedure with Stepwise Covariate Modeling (SCM) and COnditional Sampling use for Stepwise Approach (COSSAC) show similar performances on several real data examples but with a much reduced computing time. This algorithm is now implemented in Monolix and in the R package *Rsmix*.

**Study Highlights****WHAT IS THE CURRENT KNOWLEDGE ON THE TOPIC?**

Existing model-building methods for nonlinear mixed-effects models have high computational time, especially for selecting the covariate model.

**WHAT QUESTION DID THIS STUDY ADDRESS?**

The study describes the principle of the Stochastic Approximation for Model Building Algorithm (SAMBA) procedure, which allows to build a covariate, a correlation, and an error model automatically and compares it with Stepwise Covariate Modeling (SCM) and COnditional Sampling use for Stepwise Approach (COSSAC) procedures.

**WHAT DOES THIS STUDY ADD TO OUR KNOWLEDGE?**

SAMBA allows to select the best covariate model without having to fit the complete nonlinear mixed-effects model to the data for each possible covariate model. This study confirms that it is possible to obtain relevant information on the model we are looking for, even when another model is fitted to the data. This allows to drastically reduce the computation time with respect to other existing procedures

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2022 The Authors. *CPT: Pharmacometrics & Systems Pharmacology* published by Wiley Periodicals LLC on behalf of American Society for Clinical Pharmacology and Therapeutics.

while keeping the same performances. We also show that it is possible to perform correlation and error model selection in nonlinear mixed-effects models.

### **HOW MIGHT THIS CHANGE DRUG DISCOVERY, DEVELOPMENT, AND/OR THERAPEUTICS?**

This method will allow the practitioner to very quickly find a set of very good models in terms of data fitting and parsimony, even when the number of parameters or the number of covariates available is large.

## **INTRODUCTION**

Construction of a complex (nonlinear) mixed-effects model<sup>1</sup> is a challenging process which requires confirmed expertise, advanced statistical methods, and the use of sophisticated software tools, but, above all, time and patience. Indeed, the success of correctly identifying all the components of the model is far from straightforward: it is a question of finding the best structural model, determining the type of relationship between covariates and individual parameters, detecting possible correlations between random effects, or also modeling residual errors. Our goal is to accelerate and optimize this process of model building by identifying at each step how best to improve some of the model components.

The procedure for constructing a model is usually iterative: one adjusts a first model to the data, and diagnosis plots and statistical tests allow to detect possible misspecifications in the proposed model. A new model must then be proposed to correct these defects and improve the predictive abilities of the model. Most of the common approaches consist in stepwise procedures consisting in testing the addition of variable forward and their elimination backward alternatively and progressing through the choice of models using a criterion derived from the log-likelihood. A widely used approach is Stepwise Covariate Modeling (SCM),<sup>2</sup> which consists in an exhaustive search in the covariate space. Each covariate addition or deletion is tested in turn selecting models at each step leading to the best adjustment according to the objective criterion. Approaches such as Wald Approximation Method (WAM)<sup>3</sup> and COnditional Sampling use for Stepwise Approach based on Correlation tests (COSSAC)<sup>4</sup> are less computationally intensive as they use, respectively, a likelihood ratio test and a correlation test to move in the covariates space, which allows the testing of less models. All these methods are nevertheless computationally intensive as they require to re-estimate the model parameters and the likelihood many times. In particular, these methods are very sensitive to “the curse of dimensionality” when the number of covariates to test on parameters is large.

The Generalized Additive Model (GAM) method<sup>5,6</sup> is computationally appealing as it does not require as

many models fitting. Indeed, it is based on a regression on the empirical Bayes estimates (EBEs). The EBEs are the modes of the conditional distributions of the individual parameters. In other words, they are the most likely value of the individual parameters, given the estimated population parameters and the data. These estimates are known to be misleading and prone to shrinkage when data are sparse.<sup>7</sup> An efficient method which can correct the bias caused by the shrinkage of the EBEs have been recently proposed for covariate analysis.<sup>8,9</sup> In this paper, we propose to develop similar method which relies on the use of random samples from the conditional distribution of each individual parameters instead of EBEs. Indeed, the random sample of the posterior distribution has been shown to correctly control the type I error when performing tests to detect misspecifications in the model.<sup>10</sup>

As for most of the model-building procedures, the objective of Stochastic Approximation for Model Building Algorithm (SAMBA) is to find a model that minimizes some information criterion, such as Akaike information criterion (AIC), Bayesian Information Criteria (BIC), or corrected BIC (BICc).<sup>11</sup> The main principle of SAMBA is to use the results obtained with a wrong model to learn the right model. Then, SAMBA is an iterative procedure where a new model is used at each iteration of the algorithm. The values of the population parameters of the model are found by maximum likelihood estimation, and, then, the individual parameters are sampled from the conditional distribution defined under this estimated model. These simulated individual parameters combined with the observed data can now be used to select a new statistical model. It is important to underline that, as most of the iterative procedures for non-convex optimization, SAMBA does not pretend to be capable of always finding the global minimum of the used criterion, but it always allows to quickly find a very good solution.

Two contributions mainly constitute the content of this paper. First, we describe the novel algorithm called SAMBA for fast automatic model building in nonlinear mixed-effects models (section 1). Second, we benchmark its performances compared with reference methods SCM and COSSAC in real-world examples (section 2).

## METHODS

### Model description

Let  $y_i = (y_{ij}, 1 \leq j \leq n_i)$  be the vector of observations for subject  $i$ , where  $1 \leq i \leq N$ . The model that describes the observations  $y_i$  is assumed to be a parametric probabilistic model that depends on a vector of  $L$  (individual) parameters  $\psi_i = (\psi_{i1}, \dots, \psi_{iL})$ . In a population framework, the vector of parameters  $\psi_i$  is assumed to be drawn from a population distribution  $p(\psi_i)$ . Then, defining a model  $\mathcal{M}$  consists in defining a joint probability distribution for the observations  $y = (y_1, \dots, y_N)$  and for the individual parameters  $\psi = (\psi_1, \dots, \psi_N)$ . For the sake of notation simplicity, we focus on models for continuous longitudinal data. However, extension to models for discrete data and time to event data is straightforward.

Let  $y_{ij}$ , the observation obtained from subject  $i$  at time  $t_{ij}$  be described as:

$$u(y_{ij}) = u(f(t_{ij}, \psi_i)) + g(t_{ij}, \psi_i, \xi) \varepsilon_{ij}, 1 \leq i \leq N, 1 \leq j \leq n_i. \quad (1)$$

The structural model  $f$  is a fundamental component of the model because it defines the individual predictions of the observed kinetics for a given set of parameters. The residual errors ( $\varepsilon_{ij}$ ) are assumed to be standardized Gaussian random variables (mean zero and variance 1). The residual error model is represented by function  $g$  in model (1) and may depend on some additional parameter  $\xi$ . Finally, one can use the function  $u$  to transform the observations, assuming for instance that they are log-normally distributed. In the following, we will assume  $u$  to be the identity.

We assume a linear model for the individual parameters (up to some transformation  $h$ ):

$$h(\psi_i) = h(\psi_{\text{pop}}) + \beta c_i + \eta_i, 1 \leq i \leq N, \quad (2)$$

where  $\eta_i \sim \mathcal{N}(0, \Omega)$  is a vector of random effects and where  $c_i$  is a vector of individual covariates used to explain part of the variability of the  $\psi_i$ 's. The  $\psi_{\text{pop}}$  and  $\beta$  are fixed effects. The joint model of  $y$  and  $\psi$  then depends on a set of parameters  $\theta = (\psi_{\text{pop}}, \beta, \Omega, \xi)$ .

Selecting a model described by Equations 1 and 2 consists for the modeler in selecting: (i) the structural model  $f$ , (ii) the transformation of the individual parameters  $h$ , (iii) the residual error model  $g$ , (iv) the list of covariates that have an impact on individual parameters, and (v) the structure of the variance-covariance matrix of the random effects in the linear model  $\Omega$ . The selection of the two first items is problem-specific, and their selection is out of the scope of this paper. We will therefore assume, in this paper, that  $f$  and  $h$  are given. The SAMBA procedure

proposes solutions to address the selection of the three other components of the model.

### The SAMBA procedure

Automatic model building is a difficult task because it is generally not possible to fit and compare all possible models. Moreover, it is necessary to define what is the “best model” among all the possible models. A classical approach consists in searching for the model  $\mathcal{M}_*$ , that minimizes a criterion, such as the penalized likelihood<sup>12,13</sup>:

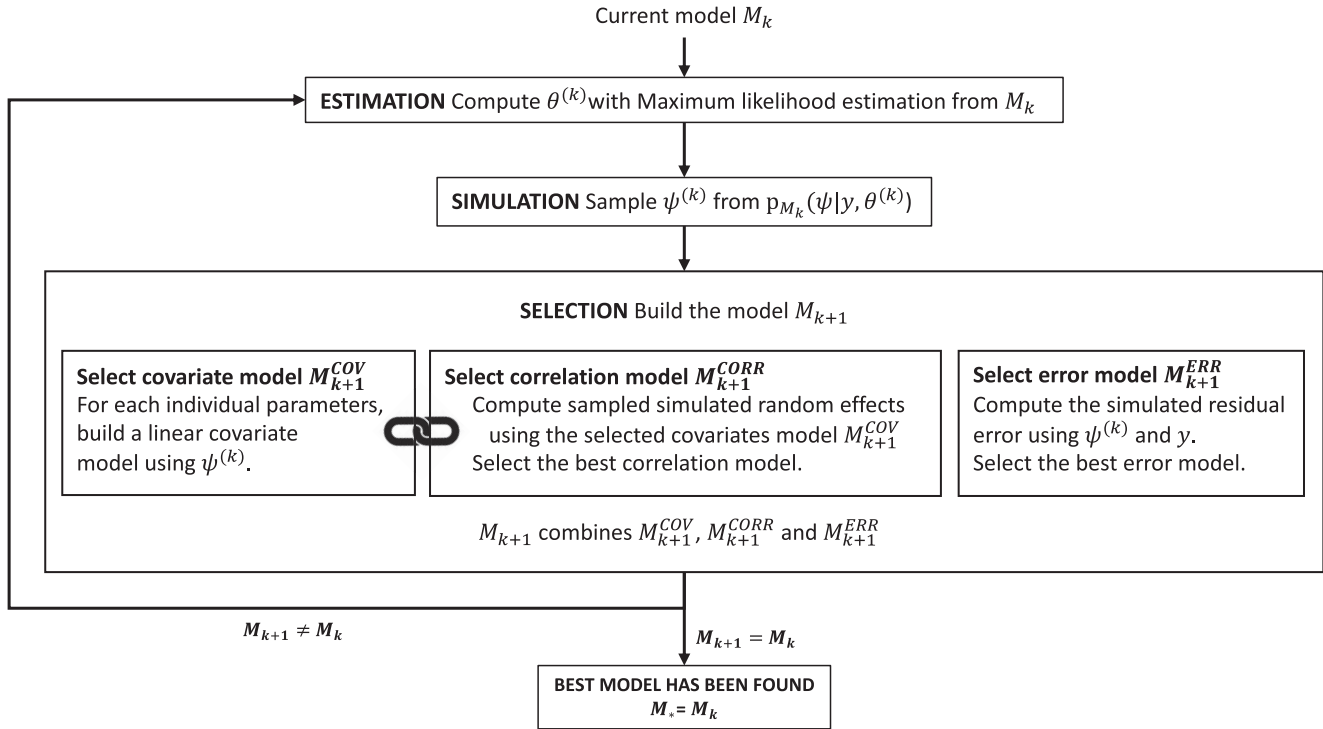
$$\mathcal{M}_* = \arg \min_{\mathcal{M}} \{ \min_{\theta} (-2 \log(\mathcal{L}_{\mathcal{M}}(\theta; y))) + \text{pen}(\mathcal{M}) \}. \quad (3)$$

The objective of this approach is to find a model that best fits the data (by minimizing  $-2LL$ ) while being as simple as possible (it is the role of  $\text{pen}(\mathcal{M})$  to favor models with few parameters). When the space of possible models is large, an exhaustive search is clearly impossible, and an efficient minimization strategy must be implemented. It is precisely for this purpose that SAMBA was developed: to obtain very quickly the “best” model  $\mathcal{M}_*$ , or a model with an objective criterion value very close to that of  $\mathcal{M}_*$ .

SAMBA is an iterative procedure alternating three steps. Assume that model  $\mathcal{M}_k$  was obtained at iteration  $k$  of the algorithm. We first compute  $\theta^{(k)}$ , the maximum likelihood estimate of  $\theta$  for model  $\mathcal{M}_k$ . We then generate a set of individual parameters  $\psi^{(k)}$  from the conditional distribution of individual parameters  $p_{\mathcal{M}_k}(\psi | y; \theta^{(k)})$ . The selection step finally consists in building a new model  $\mathcal{M}_{k+1}$  using the *complete data* ( $y; \psi^{(k)}$ ) and minimizing the complete penalized criterion:

$$\mathcal{M}_{k+1} = \arg \min_{\mathcal{M}} \{ \min_{\theta} (-2 \log(\mathcal{L}_{\mathcal{M}}(\theta; y, \psi^{(k)}))) + \text{pen}(\mathcal{M}) \}. \quad (4)$$

As already mentioned, the statistical model to be built consists of a covariate model, a correlation model, and a residual error model. Then, the selection of model  $\mathcal{M}_{k+1}$  is composed of three model selection procedures: the selection of the covariate model  $\mathcal{M}_{k+1}^{\text{COV}}$ , the selection of the correlation model  $\mathcal{M}_{k+1}^{\text{CORR}}$ , and the selection of the error model  $\mathcal{M}_{k+1}^{\text{ERR}}$ . Note that not all these components are necessarily selected: some may have been set arbitrarily because of existing knowledge. By noticing that  $\mathcal{L}_{\mathcal{M}}(\theta; y, \psi^{(k)}) = \mathcal{L}_{\mathcal{M}}(\theta | y, \psi^{(k)}) \mathcal{L}_{\mathcal{M}}(y, \psi^{(k)})$ , it appears that the problem of selecting the error model is independent from the problem of selecting the covariate and correlation models. Figure 1 provides a flowchart of the complete procedure. Let us now take a closer look at what each step of the model selection process consists of.



**FIGURE 1** Scheme of the Stochastic Approximation for Model Building Algorithm (SAMBA)

### The covariate model selection $\mathcal{M}_{k+1}^{\text{COV}}$

The sample  $\psi^{(k)}$  has been generated conditionally to the data  $y$  and the model  $\mathcal{M}_k$ . For the  $\ell$ -th parameter, we build a linear model between  $\psi_{i\ell}^{(k)}$  and covariates  $c$ , such as in Equation 2:

$$h_{\ell}(\psi_{i\ell}^{(k)}) = h_{\ell}(\psi_{\text{pop},\ell}) + \beta_{\ell} c_i + \eta_{i\ell}^{(k)}, \quad 1 \leq i \leq N, 1 \leq \ell \leq L, \quad (5)$$

with  $h_{\ell}$  the transformation associated to the  $\ell$ -th parameter and where  $\eta_{i\ell}^{(k)}$  is supposed normally distributed with mean zero and variance  $\omega_{\ell}^2$ . We define  $\theta_{\ell} = (\psi_{\text{pop},\ell}, \beta_{\ell}, \omega_{\ell}^2)$ . Best covariate model for parameter  $\ell$ , denoted  $\mathcal{M}_{k+1}^{\text{COV},\ell}$ , is selected as being the one minimizing a penalized criterion:

$$\mathcal{M}_{k+1}^{\text{COV},\ell} = \arg \min_{\mathcal{M}} \left\{ \min_{\theta_{\ell}} \left( -2 \log \left( \mathcal{L}_{\mathcal{M}} \left( \theta_{\ell}; \psi_{i\ell}^{(k)} \right) \right) \right) + \text{pen}^{\text{COV}}(\mathcal{M}) \right\}.$$

We denote  $n_{\beta}$  the number of non-null elements in  $\beta_{\ell}$  for model  $\mathcal{M}$ . The penalization depends on the criterion selected for optimization: if AIC then  $\text{pen}^{\text{COV}}(\mathcal{M}) = 2n_{\beta}$ , if BIC or BICc then  $\text{pen}^{\text{COV}}(\mathcal{M}) = \log(N)n_{\beta}$ . Equation 5 tells us that the covariate selection problem has become here a classical problem of variable selection in a linear model.<sup>14</sup> This problem is much more easily tractable than the original one. The overall best covariate model combines the best model for each parameter such that  $\mathcal{M}_{k+1}^{\text{COV}} = \left\{ \mathcal{M}_{k+1}^{\text{COV},1}, \dots, \mathcal{M}_{k+1}^{\text{COV},L} \right\}$ .

In the implemented version of package *Rsmix* (R speaks Monolix), two different strategies are implemented depending on the dimension of the selection problem. If the number  $d$  of available covariates is less than 11, an exhaustive search is performed over all the  $2^d$  possible covariate models for each parameter. Otherwise, the stepwise variable selection procedure implemented in the function *stepAIC* from package *MASS* is used. It consists of iteratively adding and removing covariates in stepwise manner to lower the objective criterion.

### The correlation model selection $\mathcal{M}_{k+1}^{\text{CORR}}$

Using the selected covariate model  $\mathcal{M}_{k+1}^{\text{COV}}$  and the sample of individual parameters  $\psi_i^{(k)}$ , it is possible to extract the vector of individual random effects  $\eta_i^{(k)} = (\eta_{i\ell}^{(k)}, \ell = 1, \dots, L)$  from Equation 5. Assuming that  $\eta_i^{(k)} \sim \mathcal{N}(0, \Omega)$  where  $\Omega$  is a block diagonal matrix, the problem of correlation model selection consists in selecting the block structure of  $\Omega$ . We then select the correlation model denoted  $\mathcal{M}_{k+1}^{\text{CORR}}$  by minimizing a penalized criterion:

$$\mathcal{M}_{k+1}^{\text{CORR}} = \arg \min_{\mathcal{M}} \left\{ \min_{\Omega} \left( -2 \log \left( \mathcal{L}_{\mathcal{M}} \left( \Omega; \eta_i^{(k)} \right) \right) \right) + \text{pen}^{\text{CORR}}(\mathcal{M}) \right\}.$$

We denote  $n_{\Omega}$  the number of non-zero elements in the upper triangular part of the matrix  $\Omega$ . The penalization

depends on the criterion selected for global optimization: if AIC then  $\text{pen}^{\text{CORR}}(\mathcal{M}) = 2n_{\Omega}$ , if BIC or BICc then  $\text{pen}^{\text{CORR}}(\mathcal{M}) = \log(N) n_{\Omega}$ .

In the implemented version of package *Rsmxl*, we limit the size of the block-structure that can be considered at each iteration. For  $\mathcal{M}_1$ , no correlation can be added and a diagonal matrix is used for  $\Omega$ ; for  $\mathcal{M}_2$  only blocks of size two are considered. At iteration  $k$  for selection of model  $\mathcal{M}_{k+1}^{\text{CORR}}$ , block size cannot be larger than  $k + 1$ , leading to no more than  $(k - 1)k/2$  non-zero covariance terms in  $\Omega$ .

### The error model selection $\mathcal{M}_{k+1}^{\text{ERR}}$

For a given set of simulated individual parameters  $(\psi_i^{(k)}, 1 \leq i \leq N)$ , the residual errors can easily be computed:

$$e_{ij}^{(k)} = y_{ij} - f(t_{ij}, \psi_i^{(k)}), 1 \leq i \leq N, 1 \leq j \leq n_i.$$

We then fit several error models with standard deviation of the form  $g(t_{ij}, \psi_i^{(k)}, \xi)$  for  $e_{ij}^{(k)}$  and select the one minimizing a penalized criterion:

$$\mathcal{M}_{k+1}^{\text{ERR}} = \text{argmin}_{\mathcal{M}} \left\{ \min_{\xi} \left( -2 \log \left( \mathcal{L}_{\mathcal{M}} \left( \xi; e_{ij}^{(k)} \right) \right) \right) + \text{pen}^{\text{ERR}}(\mathcal{M}) \right\}.$$

We denote  $n_{\xi}$  the length of  $\xi$  (i.e., the number of parameters in model  $\mathcal{M}$ ). The penalization depends on the criterion selected for global optimization: if AIC then  $\text{pen}^{\text{ERR}}(\mathcal{M}) = 2n_{\xi}$ , if BIC then  $\text{pen}^{\text{ERR}}(\mathcal{M}) = \log(N) n_{\xi}$ , and if BICc then  $\text{pen}(\mathcal{M}) = \log(n_{\text{tot}}) n_{\xi}$  where  $n_{\text{tot}}$  is the total number of observations, including below the limit of quantification data.

In the implemented version of package *Rsmxl*, five error models (provided by function `gin` Equation 1) are tested by default: constant ( $g_x(t_{ij}, \psi_i^{(k)}, \xi) = \xi$ ), proportional ( $g_x(t_{ij}, \psi_i^{(k)}, \xi) = \xi f(t_{ij}, \psi_i)$ ), combined<sub>1</sub> ( $g_x(t_{ij}, \psi_i^{(k)}, \xi) = \xi_1 + \xi_2 f(t_{ij}, \psi_i)$ ), combined<sub>2</sub> ( $g_x(t_{ij}, \psi_i^{(k)}, \xi) = \sqrt{\xi_1^2 + \xi_2^2} f(t_{ij}, \psi_i)$ ), or exponential in which a constant error model is fitted to the  $\log(y)$  using the transformation  $u = \log$  in Equation 1. Note that it is currently not possible to perform the selection on a restricted number of error models, but such a feature could be easily implemented.

### Stopping rule procedure

At each iteration  $k$  of the algorithm, we combine  $\mathcal{M}_{k+1}^{\text{COV}}$ ,  $\mathcal{M}_{k+1}^{\text{CORR}}$ , and  $\mathcal{M}_{k+1}^{\text{ERR}}$  to get the new selected model  $\mathcal{M}_{k+1}$ ,

which is passed forward on to the next estimation-simulation run. It is important to select the covariate model before the correlation model. On the other hand, the error model can be updated before or after the other two components of the model. The algorithm stops when  $\mathcal{M}_k$  is strictly identical to  $\mathcal{M}_{k+1}$  for all components and the last model is the selected one.

### Remark

In the above,  $\psi_i^{(k)}$  represents a single realization of the conditional distribution  $p_{\mathcal{M}_k}(\psi_i | y, \theta^{(k)})$  for each  $i = 1, \dots, N$ . Instead of considering only one realization of this distribution, we could use a sample of size  $R(\psi_{i\ell,r}^{(k)}, 1 \leq r \leq R)$ . If so, the linear covariate model described in Equation 5 rewrites:

$$\overline{h_{\ell}(\psi_{\ell,i}^{(k)})} = h_{\ell}(\psi_{\ell,\text{pop}}) + \beta_{\ell} c_i + \overline{\eta_{\ell,i}^{(k)}}, 1 \leq i \leq N, 1 \leq \ell \leq L,$$

where:

$$\overline{h_{\ell}(\psi_{\ell,i}^{(k)})} = \frac{1}{R} \sum_{r=1}^R h_{\ell}(\psi_{i\ell,r}^{(k)})$$

Procedures for covariate model selection and correlation model selection remains the same, but using now  $(\overline{\psi_{i\ell}^{(k)}})$  and  $(\overline{\eta_{i\ell}^{(k)}})$  at iteration  $k$ . On the other hand, the  $R$  series of residual errors  $(e_{ij,r}^{(k)})$  are used for selecting the residual error model.

## RESULTS

### Step-by-step example of the SAMBA procedure

To illustrate how SAMBA works in practice, we will describe step-by-step the complete procedure on the example of remifentanyl.<sup>15</sup> We use here the SAMBA implementation in function `buildmxx` of the R package *Rsmxl*, using the default settings.

#### The remifentanyl data

The dataset is composed of 65 healthy adults who have received remifentanyl i.v. infusion at a constant infusion rate between 1 and 8  $\mu\text{g}^{-1} \text{kg}^{-1} \text{min}^{-1}$  for 4 to 20 minutes. Time and rate of infusion are known for each individual. The pharmacokinetic (PK) data consists in the plasma concentration of remifentanyl, which is measured during and after infusion for a total of 19 to 53 observations by patients, totaling 2057 observations. A total of six

covariates are available: one qualitative covariate, the sex (SEX) and five continuous covariates: the age (AGE), the height (HT), the weight (WT), the lean body mass (LBM), and the body surface area (BSA). All the latter are normalized and log-transformed for the analysis. In the following, we adopt the notation  $\log\text{AGE} = \log(\text{AGE}/\text{AGE}_{\text{pop}})$ , where  $\text{AGE}_{\text{pop}}$  is a typical value to normalize on (e.g., the mean value of age in the population).

## The model

The PK model for i.v. infusion has a central compartment (volume V1), two peripheral compartments (volumes V2 and V3, and intercompartmental clearances Q2 and Q3), and a linear elimination (Cl). Log-normal distributions are used for the six individual parameters. The  $2^6 = 64$  possible covariate models will be considered for each of the six individual parameters. Note that if we had to test all possible models, we would have had to test  $64^6$  combinations, which would have made the problem intractable.

## SAMBA iterations

We start the SAMBA procedure with a model  $\mathcal{M}_0$  without any covariate on all parameters, with no correlation between random effects and the so-called combined<sub>1</sub> error model. Figure 2 illustrates the selection steps on this specific example. One can notice that the BICc, which has been chosen as target criterion, decreases from 7186 for  $\mathcal{M}_0$  to 6985 for  $\mathcal{M}_1$ , 6957 for  $\mathcal{M}_2$ , and 6903 for  $\mathcal{M}_3$ , which is finally selected as the best model for this example.

- **Run 0 (BICc = 7185.8) + Iteration 1:** Model  $\mathcal{M}_0$  is fitted to data and individual parameters are sampled conditionally on the data and this model. Each of the 64 possible linear covariate models is fitted to each individual parameters and the one with lowest BICc is selected. Let us take the example of Cl: the three best models include (1) an effect of logAGE and logWT (BICc = -55.0), (2) an effect of logAGE and logLBM (BICc = -56.1), and (3) an effect of logAGE and logBSA (BICc = -57.5). The latter is chosen as the best model for parameter Cl as it provides the lowest BICc ( $\mathcal{M}_1^{\text{COV},Cl}$ ). Altogether, for all parameters, the best covariate model ( $\mathcal{M}_1^{\text{COV}}$ ) includes logAGE on all parameters, logBSA on Cl, and logLBM on V1 and V2. No correlation is added to the model because no correlation is allowed at first iteration. Then,  $\mathcal{M}_1^{\text{CORR}}$  is a diagonal variance-covariance matrix for the random effects. Among the tested error models, the

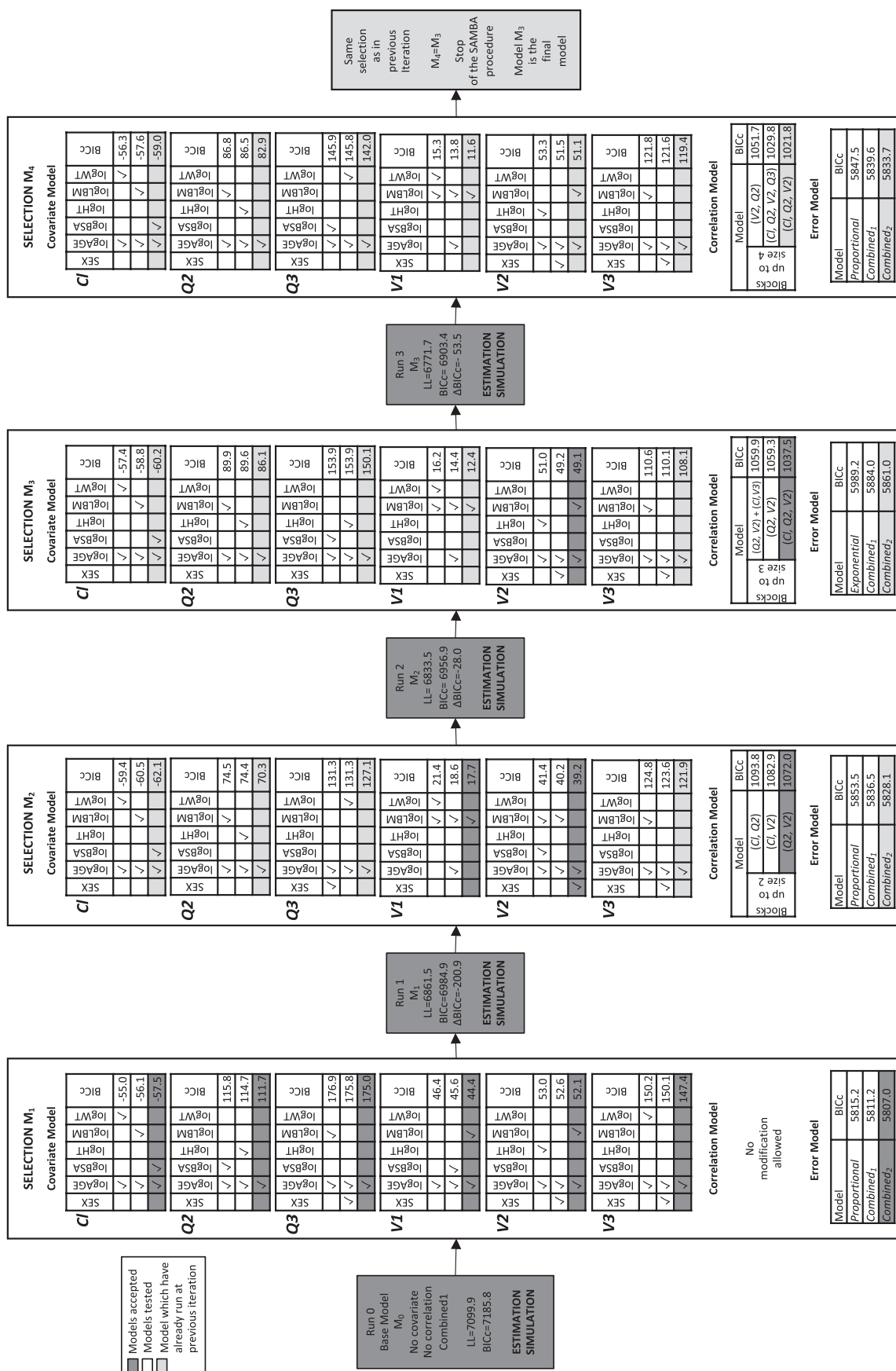
three best ones are proportional (BICc = 5815.2), combined<sub>1</sub> (BICc = 5811.2), and combined<sub>2</sub> (BICc = 5807.0), which is selected for  $\mathcal{M}_1^{\text{ERR}}$ . These covariate, correlation, and error models are then passed on to run 1:

$$\mathcal{M}_1 = \{\{\mathcal{M}_1^{\text{COV},Cl}, \mathcal{M}_1^{\text{COV},Q2}, \mathcal{M}_1^{\text{COV},Q3}, \mathcal{M}_1^{\text{COV},V1}, \mathcal{M}_1^{\text{COV},V2}, \mathcal{M}_1^{\text{COV},V3}, \mathcal{M}_1^{\text{CORR}}, \mathcal{M}_1^{\text{ERR}}\}\}.$$

- **Run 1 (BICc = 6984.9) + Iteration 2:** Model  $\mathcal{M}_1$  is fitted to the data and individual parameters are sampled. Again, the three best model for each covariate are provided. The best covariate model includes logAGE on all parameters except V1, logBSA on Cl, logLBM on V1, and SEX on V2 ( $\mathcal{M}_2^{\text{COV}}$ ). Block-structured correlation with blocks up to size 2 are compared (i.e., up to one correlation term). The best three models are with a correlation between parameters Cl and V2 (BICc = 1082.9), between parameters Cl and Q2 (BICc = 1093.8), and between parameters V2 and Q2 (BICc = 1072.0). The latter correlation model is selected for  $\mathcal{M}_2^{\text{CORR}}$ . Residual error model combined<sub>2</sub> remains the best one ( $\mathcal{M}_2^{\text{ERR}}$ ). These covariate, correlation, and error models are then passed on to run 2.
- **Run 2 (BICc = 6956.9) + Iteration 3:** Model  $\mathcal{M}_2$  is fitted to data and individual parameters are sampled. The best covariate model includes logAGE on all parameters except V1, logBSA on Cl, and logLBM on V1 and V2 ( $\mathcal{M}_3^{\text{COV}}$ ). Block-structured correlation with blocks up to size 3 are compared (i.e., up to three correlation terms), a correlation block is selected between Cl, Q2, and V2 ( $\mathcal{M}_3^{\text{CORR}}$ ). Residual error model combined<sub>2</sub> remains the best one ( $\mathcal{M}_3^{\text{ERR}}$ ). These covariate, correlation, and error models are then passed on to run 3.
- **Run 3 (BICc = 6903.4) + Iteration 4:** Model  $\mathcal{M}_3$  is fitted to data and individual parameters are sampled. Of note, regarding the correlation model selection, block-structured correlation with blocks up to size 4 are compared (i.e., up to six correlation terms). During this iteration, the same model as the one in the previous iteration is selected ( $\mathcal{M}_4 = \mathcal{M}_3$ ) resulting in the stopping of the procedure. Model  $\mathcal{M}_3$  is therefore the final model selected with the SAMBA procedure.

## Converging toward a global optimal model

Even if the selected criterion decreases at each iteration, there is no guarantee that SAMBA converges toward a global minimum of this criterion. The quality and the robustness of the convergence of SAMBA can then be assessed by running SAMBA several times from different starting models. In particular, a good practice is to: (1) launch SAMBA from several initial models, (2) compare the best models found (if there is not only one) in terms of objective criterion (e.g., BICc), and (3) make



**FIGURE 2** Step-by-step Stochastic Approximation for Model Building Algorithm (SAMBA) procedure on the remifentanyl example with six covariates (SEX, logAGE, logBSA, logHT, logLBM, and logWT) and six model parameters (C1, Q2, Q3, V1, V2 and V3). For each selection (covariate, correlation, and error model), the three best models in term of corrected Bayesian Information Criteria (BICc) are displayed. Non selected models are in darker grey, and models which have been already accepted at previous run are in lighter grey

a thorough analysis and interpretation of the nearby models in order to choose the most relevant one for a given application. Regarding the choice of the starting model, similarly to the Expectation Maximization and Stochastic Approximation Expectation Maximization algorithms, there is no optimal choice.<sup>16,17</sup> We recommend to test in priority the following three starting models: (1) an empty model, (2) (when possible) a complete model, and (3) a model (or models) that make sense for the biological application. Note that this robustness assessment is standard for all non-convex optimization algorithms and should also be performed for SCM and COSSAC in routine.

### Performances on real examples, and comparison with the SCM and COSSAC procedures

To assess the performances of the SAMBA procedure compared to SCM and COSSAC procedures, we replicate the illustration provided in ref. 4. We applied the three routines to a collection of 10 representative datasets, including PKs, pharmacodynamics, and disease models. Of note, the SCM method for variable selection used here is exactly the same as the one implemented in PsN (Pearl Speaks NONMEM), differences lie in the algorithms used to estimate the parameters of a model and to calculate the likelihood. We restricted the SAMBA procedure to the covariate model selection as correlation and error model selection are not implemented in COSSAC and SCM. The results can be found in Table 1.

Because the datasets are real data illustrations, there is no “true” model. It is only possible to compare them in terms of BIC. Of 10 examples, the same best model was proposed by the three procedures in four examples. In two examples, the best model selected by SAMBA was better in terms of BICc than with SCM and COSSAC (Theophylline Ext. Rel. and Warfarin PK/PD). In three other examples, the model with the lowest BICc was not selected by SAMBA. However, the difference in BICc was, respectively, smaller than six in comparison with the SCM procedure and 4.2 in comparison with the COSSAC procedure. We insist on the fact that a difference in BICc does not necessarily have any biological meaning. This is an arbitrary criterion that allows to quantify the goodness of fit with respect to the sparsity of the model chosen. We thus argue that the three procedures lead to rather similar models, which all constitute very good starting points for the modeler to build a model based on biological hypothesis. Finally, in only one example discussed below, the difference in BICc was larger than 10 points of BICc both compared with the SCM and COSSAC procedures.

Regarding the cholesterol dataset, we again ran the SAMBA procedure starting from a full model in which all covariates are supposed to have an effect on all parameters. The new model selected by SAMBA is the full model with an effect of logAGE on (*Chol0*, *slope*) and SEX on (*Chol0*, *slope*) is much closer in term of BICc than the one selected starting from an empty model ( $\Delta \text{BICc} = -2$ ). We can finally notice with this example that it is sometimes possible to improve the convergence of SAMBA by improving the convergence of SAEM. Indeed, using 10 Markov chains instead of only one, SAMBA also finds the model selected by SCM and COSSAC. Finding the optimal settings that minimizes computation time while maximizing the probability of finding the best model is an extremely difficult problem that remains open. We can claim that the default settings used in Rsmx and Monolix give very good results in most cases, but not in all cases with absolute certainty.

In terms of computational effort, it is important to note that the SAMBA procedure completes the model-building process in much less runs, hence much less CPU time than SCM and COSSAC. In the considered problems, the number of runs and the CPU computation time are equivalent because the other computation times are negligible in the order of a few seconds. Actually, the computation times are six to 149 smaller than for SCM and two to 11 times smaller than for COSSAC. Note that the number of evaluations required by SAMBA is always lower or equal to the number of evaluations performed by COSSAC and SCM.

### Simulation study

#### Data generation and analysis

We simulated data from a one-compartment PK model. The model has three population parameters  $ka_{\text{pop}} = 1$ ,  $V_{\text{pop}} = 10$  and  $Cl_{\text{pop}} = 2$ . All individual parameters are log-normally distributed around the population parameters ( $\omega_{ka} = 0.2$ ,  $\omega_V = 0.3$  and  $\omega_{Cl} = 0.3$ ). We simulated five individual covariates ( $C_1, C_2, C_3, C_4, C_5$ ) from standard normal distributions. The covariate model is such that there only exists linear relationships between  $\log(V)$  and  $C_1$  ( $\beta_{V,1} = 0.2$ ),  $\log(Cl)$  and  $C_1$  ( $\beta_{Cl,1} = -0.2$ ), and  $\log(Cl)$  and  $C_2$  ( $\beta_{Cl,2} = 0.3$ ). The correlation model is such that there exists a linear correlation between  $\eta_V$  and  $\eta_{Cl}$  ( $\rho_{V,Cl} = 0.6$ ). Finally, the error model is a combined<sub>2</sub> model with  $a = 2$  and  $b = 0.1$ . A clinical trial could then be simulated by generating PK data from this model for 100 individuals and 11 timepoints (0.25, 0.5, 1, 2, 5, 8, 12, 16, 20, 24, and 30). In order to evaluate the properties of SAMBA by Monte-Carlo, we simulated 100 replicates of the same trial and built the model for each replicate using SAMBA as implemented in Rsmx and Monolix for minimizing BICc.



**TABLE 1** Comparison of the SAMBA procedure with the SCM and COSSAC procedure on 10 representative datasets

Dataset	Characteristics	SCM		COSSAC		SAMBA		ΔBICc	
		#Runs <sup>b</sup>	Final Model <sup>a</sup>	#Runs <sup>b</sup>	Final Model 1	#Runs <sup>b</sup>	Final Model <sup>a</sup>	SAMBA-SCM	SAMBA-COSSAC
Warfarin	32 ind. - 247 obs. 4 param. - 3 cov. 4 re - 1 outcome	44	logWt - V, Cl logAge - C	4	Identical	2	Identical	0	0
Remifentanyl	65 ind. - 1992 obs.	295	logLBM - V1	13	logLBM - V1, V2	4	logLBM - V1	0.8	0.5
Linear PK	6 param. - 6 cov.		logAGE - Cl, Q2, Q3, V2, V3		logAGE - Cl, Q2, V2, V3		logAGE - Cl, Q2, Q3, V2, V3		
	4 re - 1 outcome		logBSA - Cl logHT - V2		logBSA - Cl		logBSA - Cl		
Theophylline	12 ind. - 20 obs. 3 param. - 2 cov. 4 re - 1 outcome	12	logWEIGHT - ka	4	Identical	2	Identical	0	0
Quinidine	136 ind. - 361 obs. 3 param. - 2 cov. 3 re - 1 outcome	22	none	11	Identical	1	Identical	0	0
Tobramycin	97 ind. - 322 obs. 3 param. - 2 cov. 2 re - 1 outcome	22	logCLCR - Cl	6	logCLCR - Cl	2	logCLCR - Cl	4.2	4.2
			logWT - V		logWT - V		logWT - CI		
Theophylline	18 ind. - 362 obs. 7 param. - 3 cov. 7 re - 1 outcome	98	logWT - Tlag1, V	8	logWT - Tlag1	6	logWT - F, V logAGE - F	-11.7	-27
					logAGE - ka2		logHT ka1, ka2, Tlag1, diffTlag2		
Warfarin	32 ind. - 247+232 obs.	92	logWT - Cl	10	logWT - Cl	2	logWT - Cl, V	-1.4	-1.4
PK/PD	8 param. - 3 cov. 8 re - 2 outcomes						logAGE - CI, R0		
Cholesterol	200 ind. - 1044 obs.	12	logAGE - Chol0, slope	5	logAGE - Chol0, slope	2	logAGE - Chol0	13.5	13.5
Disease	2 param. - 2 cov.		SEX - slope		SEX - slope				
Progression	2 re - 1 outcome								
Alzheimer	896 ind. - 3707 obs.	73	APOE - alpha, p0	8	APOE - alpha, p0	2	APOE - alpha, p0	6	1.5

(Continues)

TABLE 1 (Continued)

Dataset	SCM		COSSAC		SAMBA		ABICc			
	Characteristics	#Runs <sup>b</sup>	Final Model <sup>a</sup>	#Runs <sup>b</sup>	Final Model 1	#Runs <sup>b</sup>	Final Model <sup>a</sup>	SAMBA-COSSAC	SAMBA-SCM	SAMBA-COSSAC
Sparse PK	2 param. - 7 cov.		logAGE - p0, alpha		logAGE - p0, alpha		logAGE - p0			
	2 re - 1 outcome		logBMI - alpha		logBMI - alpha		logWT - p0			
Tranexamic	166 ind. - 817 obs.	298	logWT - p0	12	logWT - p0	2	Identical	0		0
PK	4 param. - 10 cov.		GROUP - Cl, V2		Identical		Identical			
	4 re - 1 outcome		logBMI - Cl							
			logCOCK - Cl							
			logLBW - Q							
			logWeight - V2							

<sup>a</sup>Differences of variable selection between different methods are highlighted in bold.

<sup>b</sup>The number of runs is defined as the number of time the estimation and the simulation steps are performed (which is the most time-consuming).

TABLE 2 Performance of the SAMBA algorithm for the selection of the covariate model in a simulation study using a one-compartment PK model

Covariates	<i>Rsmix</i>			Monolix		
	ka	V	Cl	ka	V	Cl
C <sub>1</sub>	2	100	100	2	100	100
C <sub>2</sub>	0	1	100	0	1	100
C <sub>3</sub>	1	2	1	2	2	1
C <sub>4</sub>	0	3	4	0	3	4
C <sub>5</sub>	0	1	1	1	2	1

One hundred datasets of 100 individuals with 11 observations each have been generated. True model  $\mathcal{M}^*$  includes an effect of C<sub>1</sub> on V and Cl and an effect of C<sub>2</sub> on Cl. The percentages of times (over 100 replicates) each covariate-parameter relationship is selected in the final model are displayed. Implementation of SAMBA in *Rsmix* and Monolix are compared.

Abbreviations: Cl, linear elimination; ka, absorption rate constant; PK, pharmacokinetic; SAMBA, Stochastic Approximation for Model Building Algorithm; V, volume.

The initial model did not include any covariate-parameter relationship and any correlation between random effect. The initial residual error model was a combined<sub>1</sub> model. The R code used for this Monte-Carlo study is available as Supplementary Material.

## Performances

Table 2 summarizes the results obtained for the covariate model selection. On the one hand, we can see that, for this particular example, SAMBA finds the three existing covariate-parameter relationships in 100% of the cases. On the other hand, very few spurious relationships are detected (less than 2%). Importantly, in all cases for which the final covariate model included more covariates than the true model  $\mathcal{M}^*$ , the BICc of the selected model was lower than that of  $\mathcal{M}^*$  (the differences ranging from 3 to 14.7 with *Rsmix* and from 2.4 to 14.6 for Monolix). In other words, SAMBA always finds a covariate model as good or better than  $\mathcal{M}^*$  in terms of BICc. Regarding the selection of the correlation model, the correct model was selected for all the replicates. Finally, the correct error model was selected in 86% of the times with *Rsmix* and 85% of the times with Monolix. Note that all the wrong selected error models were all combined<sub>1</sub> model (instead of combined<sub>2</sub>) with a slightly larger BICc most of the time. Actually, these two models are quite similar and difficult to distinguish on the basis of a criterion like BICc. SAMBA then may get stuck in a local minimum in such a situation. Finally, and importantly, the final selected models obtained with *Rsmix* and Monolix are different in only 6% of cases. These small differences are due to small

differences in the implementation of the algorithm (see the Discussion section for more details).

## DISCUSSION

This paper presents a novel model-building procedure which offers covariate, correlation, and error model selection. It is fast as it requires only a limited number of runs of population parameter estimation and simulation compared to SCM and COSSAC. It allows to explore the space of models rapidly and provides to the modeler a very good model in term of the selection criterion. However, we insist on the fact that this procedure does not aim at replacing model-building based on biological knowledge, which is, in essence, the strength of mechanistic modeling. Thus, it should not be blindly used and the best—potentially few best—models should be interpreted and compared.

SAMBA is an efficient algorithm for minimizing an objective function. In this paper, we do not aim at evaluating the quality of the criterion used for model selection.<sup>18</sup> What is of interest here is the convergence of SAMBA. As it is also the case for SCM and COSSAC, SAMBA may not converge to the global minimum. This is particularly the case when the amount of data is too small compared to the complexity of the model to build. This phenomenon will be particularly critical when the number of covariates is high and/or when these are highly correlated. We then strongly encourage the user to build strategies to assess the robustness of the results. Extensions of the proposed algorithm are possible but are outside the scope of this paper and constitute a possible new research direction.

When there is a large number of available covariates, COSSAC and mainly SCM often fail in finding the best model in a reasonable time. In this case, SAMBA represents a particularly appealing approach because the covariate model selection is based on a stepwise variable selection procedure for linear models, which is known to handle high-dimension problems. Although stepwise AIC/BIC are designed to obtain a sparse estimator that works well on the training set, other methods, such as the lasso,<sup>19</sup> where the penalty is chosen with cross validation, is designed to obtain the sparse linear model that minimize the prediction error. A lasso type approach<sup>20</sup> can sometimes present better performances than an approach based on an information criterion, such as AIC or BIC, in particular when the number of covariates is very high. However, it should be noted that the choice of the penalty parameter by cross-validation can be complicated to implement and require a large number of runs. This type of method could be alternatively implemented in the covariate selection procedure and compared

in further works. Note finally that it would be interesting to study the behavior of SAMBA using the EBEs (corrected as proposed in ref. 8,9), rather than the individual simulated parameters, to build the covariate model.

The SAMBA procedure is implemented the R Package *Rsmx* in the function *buildmrx*.<sup>21</sup> Minimal required input is a Monolix project used as initial model. Additional arguments can be used to enable specific features (all not listed): select the components of the model to optimize among the covariate, correlation, and error model, restrict the number of parameters or covariates to use, select a specific objective criterion, etc. *Rsmx* is on CRAN and the R code can be modified to investigate any of the alternative implementations mentioned above for a specific problem. Note that the execution of *Rsmx* requires the Monolix software, because it is only an algorithm combining tasks implemented in Monolix. The R codes allowing to replicate the analyses of this paper are available in the Supplementary Material. All the illustration datasets can be downloaded from the Supporting Information Appendix S2 of ref. 4.

Finally, the SAMBA procedure is also implemented in the Monolix-GUI software starting from version 2019. Implementation is similar to the one in *Rsmx* with two noteworthy differences. First, for the selection of covariates, a stepwise procedure is used even if the number of covariates  $d$  is small. Second, compiling differences exist between C++ and R. The full SAMBA procedure is available in the model-building perspective, under a task called automatic statistical model building method. A single iteration of the SAMBA procedure is also proposed in the section Proposal in the tab Results after running a single estimation and simulation step for a model in Monolix.<sup>22</sup>

## CONFLICT OF INTEREST

Marc Lavielle is chief scientist of Lixoft, the company that develops and distributes the Monolix Suite. The other author declared no competing interests for this work.

## AUTHOR CONTRIBUTION

M.P. and M.L. wrote the manuscript, designed the research, performed the research, and analyzed the data. M.L. contributed new reagents/analytical tools.

## REFERENCES

1. Lavielle M. *Mixed-effects Models for the Population Approach: Models, Tasks, Methods and Tools*. CRC Press; 2014.
2. Jonsson EN, Karlsson MO. Automated covariate model building within NONMEM. *Pharm Res*. 1998;15(9):1463-1468.
3. Kowalski KG, Huttmacher MM. Efficient screening of covariates in population models using Wald's approximation to the likelihood ratio test. *J Pharmacokinetics Pharmacodyn*. 2001;28(3):253-275.
4. Ayral G, Si Abdallah J-F, Magnard C, Chauvin J. A novel method based on unbiased correlations tests for covariate selection in

- nonlinear mixed-effects models—the COSSAC approach. *CPT Pharmacometrics Syst Pharmacol.* 2021;10(4):318-329.
5. Hastie T, Tibshirani R. Generalized additive models: some applications. *J Am Stat Assoc.* 1987;82(398):371-386.
  6. Mandema JW, Verotta D, Sheiner LB. Building population pharmacokinetic/pharmacodynamic models. I. models for covariate effects. *J Pharmacokinetic Biopharm.* 1992;20(5):511-528.
  7. Nguyen T, Mouksassi M-S, Holford N, et al. Model evaluation of continuous data pharmacometric models: metrics and graphics. *CPT: Pharmacometrics Syst Pharmacol.* 2017;6(2):87-109.
  8. Yuan M, Xu XS, Yang Y, et al. A quick and accurate method for the estimation of covariate effects based on empirical bayes estimates in mixed-effects modeling: correction of bias due to shrinkage. *Stat Methods Med Res.* 2019;28(12):3568-3578.
  9. Yuan M, Zhu Z, Yang Y, et al. Efficient algorithms for covariate analysis with dynamic data using nonlinear mixed-effects model. *Stat Methods Med Res.* 2021;30(1):233-243.
  10. Lavielle M, Ribba B. Enhanced method for diagnosing pharmacometric models: random sampling from conditional distributions. *Pharm Res.* 2016;33(12):2979-2988.
  11. Delattre M, Lavielle M, Poursat M-A. A note on BIC in mixed-effects models. *Elect J Stat.* 2014;8(1):456-475.
  12. Green PJ. Penalized likelihood. *Encyclopedia Stat Sci.* 1998;2:578-586.
  13. James G, Witten D, Hastie T, Tibshirani R. *An Introduction to Statistical Learning, volume 112.* Springer; 2013.
  14. George EI. The variable selection problem. *J Am Stat Assoc.* 2000;95(452):1304-1308.
  15. Minto CF, Schnider TW, Egan TD, et al. Influence of age and gender on the pharmacokinetics and pharmacodynamics of remifentanyl: I. model development. *J Am Soc Anesthesiol.* 1997;86(1):10-23.
  16. Baudry J-P, Celeux G. EM for mixtures. *Stat Comp.* 2015;25(4):713-726.
  17. Biernacki C, Celeux G, Govaert G. Choosing starting values for the EM algorithm for getting the highest likelihood in multivariate gaussian mixture models. *Comput Stat Data Anal.* 2003;41(3-4):561-575.
  18. Buatois S, Ueckert S, Frey N, Retout S & Mentré F. Comparison of model averaging and model selection in dose finding trials analyzed by nonlinear mixed effect models. *The AAPS J.* 2018;20(3):1-9.
  19. Tibshirani R. Regression shrinkage and selection via the lasso. *J Roy Stat Soc: Ser B (Methodol).* 1996;58(1):267-288.
  20. Hastie T, Tibshirani R, Tibshirani R. Best subset, forward stepwise or lasso? analysis and recommendations based on extensive comparisons. *Stat Sci.* 2020;35(4):579-592.
  21. Lavielle M. Rsmlx: R Speaks 'Monolix'. 2021. <http://rsmlx.webpobox.org>. R package version 3.0.3.
  22. Monolix Online Documentation - Proposal Tab Description. <https://monolix.lixoft.com/tasks/proposal/>. Accessed September 15, 2021.

### SUPPORTING INFORMATION

Additional supporting information may be found in the online version of the article at the publisher's website.

**How to cite this article:** Prague M, Lavielle M. SAMBA: A novel method for fast automatic model building in nonlinear mixed-effects models. *CPT Pharmacometrics Syst Pharmacol.* 2022;11:161-172. doi:[10.1002/psp4.12742](https://doi.org/10.1002/psp4.12742)