



IN PARTNERSHIP WITH:
Ecole Polytechnique

Activity Report 2018

Project-Team XPOP

Statistical modelling for life sciences

IN COLLABORATION WITH: Centre de Mathématiques Appliquées (CMAP)

RESEARCH CENTER
Saclay - Île-de-France

THEME
Modeling and Control for Life Sciences

Table of contents

1. Team, Visitors, External Collaborators	1
2. Overall Objectives	2
2.1. Developing sound, useful and usable methods	2
2.2. Combining numerical, statistical and stochastic components of a model	2
2.3. Developing future standards	3
3. Research Program	3
3.1. Scientific positioning	3
3.2. The mixed-effects models	4
3.3. Computational Statistical Methods	4
3.4. Markov Chain Monte Carlo algorithms	5
3.5. Parameter estimation	5
3.6. Model building	6
3.7. Model evaluation	7
3.8. Missing data	7
4. Application Domains	8
4.1. Precision medicine and pharmacogenomics	8
4.2. Oncology	9
4.3. Hemodialysis	9
4.4. Intracellular processes	10
4.5. Population pharmacometrics	10
4.6. Mass spectrometry	11
5. Highlights of the Year	11
6. New Software and Platforms	11
6.1. mlxR	11
6.2. Rsmlx	12
6.3. SPIX	12
7. New Results	12
7.1. Normalizing constants of log-concave densities	12
7.2. The Tamed Unadjusted Langevin Algorithm	12
7.3. Development and performance of npde for the evaluation of time-to-event models	13
7.4. Low-rank Interaction with Sparse Additive Effects Model for Large Data Frames	13
7.5. Diffusion approximations and control variates for MCMC	13
7.6. Density estimation for random walks in random environment	13
7.7. Imputation of mixed data with multilevel singular value decomposition	14
7.8. Logistic Regression with Missing Covariates – Parameter Estimation, Model Selection and Prediction	14
7.9. A fast Stochastic Approximation of the EM algorithm for nonlinear mixed effects models	14
7.10. Incomplete graphical model inference via latent tree aggregation	15
8. Bilateral Contracts and Grants with Industry	15
9. Partnerships and Cooperations	15
9.1. National Initiatives	15
9.1.1. ANR	15
9.1.2. Institut National du Cancer (INCa)	15
9.2. International Initiatives	15
9.3. International Research Visitors	16
10. Dissemination	16
10.1. Promoting Scientific Activities	16
10.1.1. Scientific Events Organisation	16
10.1.2. Scientific Events Selection	16

10.1.3. Journal	16
10.1.4. Invited Talks	16
10.1.5. Leadership within the Scientific Community	17
10.1.6. Scientific Expertise	17
10.1.7. Research administration	17
10.2. Teaching - Supervision - Juries	17
10.2.1. Teaching	17
10.2.2. Supervision	17
10.3. Popularization	17
10.3.1. Internal or external Inria responsibilities	17
10.3.2. Creation of media or tools for science outreach	18
11. Bibliography	18

Project-Team XPOP

Creation of the Team: 2016 January 01, updated into Project-Team: 2017 July 01

Keywords:

Computer Science and Digital Science:

- A3.1.1. - Modeling, representation
- A3.2.3. - Inference
- A3.3. - Data and knowledge analysis
 - A3.3.1. - On-line analytical processing
 - A3.3.2. - Data mining
 - A3.3.3. - Big data analysis
- A3.4.1. - Supervised learning
- A3.4.2. - Unsupervised learning
- A3.4.4. - Optimization and learning
- A3.4.5. - Bayesian methods
- A3.4.6. - Neural networks
- A3.4.7. - Kernel methods
- A3.4.8. - Deep learning
- A5.9.2. - Estimation, modeling
- A6.1.1. - Continuous Modeling (PDE, ODE)
- A6.2.2. - Numerical probability
- A6.2.3. - Probabilistic methods
- A6.2.4. - Statistical methods
- A6.3.3. - Data processing
- A6.3.5. - Uncertainty Quantification

Other Research Topics and Application Domains:

- B1.1.4. - Genetics and genomics
- B1.1.7. - Bioinformatics
- B1.1.10. - Systems and synthetic biology
- B2.2.3. - Cancer
- B2.2.4. - Infectious diseases, Virology
- B2.4.1. - Pharmacokinetics and dynamics
- B9.1.1. - E-learning, MOOC

1. Team, Visitors, External Collaborators

Research Scientists

Marc Lavielle [Team leader, Inria, Senior Researcher]
Laetitia Le [APHP - Ecole polytechnique, Researcher, from Nov 2018]

Faculty Members

Julie Josse [Ecole polytechnique, Associate Professor]
Erwan Le Pennec [Ecole Polytechnique, Associate Professor, HDR]
Eric Moulines [Ecole polytechnique, Professor]

Post-Doctoral Fellow

Tom Rohmer [Inria, from Sep 2018]

PhD Students

Nicolas Brosse [Ecole polytechnique]

Wei Jiang [Ecole polytechnique]

Mohammed Karimi [Inria]

Genevieve Robin [Ecole polytechnique]

Marine Zulian [Dassault Systemes]

Technical staff

Yao Xu [Inria]

Intern

Duc Dung Tran [Inria, from Apr 2018 until Jul 2018]

Administrative Assistants

Hanadi Dib [Inria, until Oct 2018]

Ines Dumontier [Inria, from Nov 2018]

Hélène Kutniak [Inria, from Oct 2018]

Visiting Scientist

Ricardo Rios [Universidad Central de Venezuela, Sep 2018]

2. Overall Objectives

2.1. Developing sound, useful and usable methods

The main objective of XPOP is to develop new sound and rigorous methods for statistical modeling in the field of biology and life sciences. These methods for modeling include statistical methods of estimation, model diagnostics, model building and model selection as well as methods for numerical models (systems of ordinary and partial differential equations). Historically, the key area where these methods have been used is population pharmacokinetics. However, the framework is currently being extended to sophisticated numerical models in the contexts of viral dynamics, glucose-insulin processes, tumor growth, precision medicine, spectrometry, intracellular processes, etc.

Furthermore, an important aim of XPOP is to transfer the methods developed into software packages so that they can be used in everyday practice.

2.2. Combining numerical, statistical and stochastic components of a model

Mathematical models that characterize complex biological phenomena are defined by systems of ordinary differential equations when dealing with dynamical systems that evolve with respect to time, or by partial differential equations when there is a spatial component in the model. Also, it is sometimes useful to integrate a stochastic aspect into the dynamical system in order to model stochastic intra-individual variability.

In order to use such methods, we must deal with complex numerical difficulties, generally related to resolving the systems of differential equations. Furthermore, to be able to check the quality of a model (i.e. its descriptive and predictive performances), we require data. The statistical aspect of the model is thus critical in how it takes into account different sources of variability and uncertainty, especially when data come from several individuals and we are interested in characterizing the inter-subject variability. Here, the tools of reference are mixed-effects models.

Confronted with such complex modeling problems, one of the goals of XPOP is to show the importance of combining numerical, statistical and stochastic approaches.

2.3. Developing future standards

Linear mixed-effects models have been well-used in statistics for a long time. They are a classical approach, essentially relying on matrix calculations in Gaussian models. Whereas a solid theoretical base has been developed for such models, *nonlinear* mixed-effects models (NLMEM) have received much less attention in the statistics community, even though they have been applied to many domains of interest. It has thus been the users of these models, such as pharmacometricians, who have taken them and developed methods, without really looking to develop a clean theoretical framework or understand the mathematical properties of the methods. This is why a standard estimation method in NLMEM is to linearize the model, and few people have been interested in understanding the properties of estimators obtained in this way.

Statisticians and pharmacometricians frequently realize the need to create bridges between these two communities. We are entirely convinced that this requires the development of new standards for population modeling that can be widely accepted by these various communities. These standards include the language used for encoding a model, the approach for representing a model and the methods for using it:

- **The approach.** Our approach consists in seeing a model as hierarchical, represented by a joint probability distribution. This joint distribution can be decomposed into a product of conditional distributions, each associated with a submodel (model for observations, individual parameters, etc.). Tasks required of the modeler are thus related to these probability distributions.
- **The methods.** Many tests have shown that algorithms implemented in MONOLIX are the most reliable, all the while being extremely fast. In fact, these algorithms are precisely described and published in well known statistical journals. In particular, the SAEM algorithm, used for calculating the maximum likelihood estimation of population parameters, has shown its worth in numerous situations. Its mathematical convergence has also been proven under quite general hypotheses.
- **The language.** Mlxtran is used by MONOLIX and other modeling tools and is today by far the most advanced language for representing models. Initially developed for representing pharmacometric models, its syntax also allows it to easily code dynamical systems defined by a system of ODEs, and statistical models involving continuous, discrete and survival variables. This flexibility is a true advantage both for numerical modelers and statisticians.

3. Research Program

3.1. Scientific positioning

"Interfaces" is the defining characteristic of XPOP:

The interface between statistics, probability and numerical methods. Mathematical modelling of complex biological phenomena require to combine numerical, stochastic and statistical approaches. The CMAP is therefore the right place to be for positioning the team at the interface between several mathematical disciplines.

The interface between mathematics and the life sciences. The goal of XPOP is to bring the right answers to the right questions. These answers are mathematical tools (statistics, numerical methods, etc.), whereas the questions come from the life sciences (pharmacology, medicine, biology, etc.). This is why the point of XPOP is not to take part in mathematical projects only, but also pluridisciplinary ones.

The interface between mathematics and software development. The development of new methods is the main activity of XPOP. However, new methods are only useful if they end up being implemented in a software tool. On one hand, a strong partnership with Lixoft (the spin-off company who continue developing MONOLIX) allows us to maintaining this positioning. On the other hand, several members of the team are very active in the R community and develop widely used packages.

3.2. The mixed-effects models

Mixed-effects models are statistical models with both fixed effects and random effects. They are well-adapted to situations where repeated measurements are made on the same individual/statistical unit.

Consider first a single subject i of the population. Let $y_i = (y_{ij}, 1 \leq j \leq n_i)$ be the vector of observations for this subject. The model that describes the observations y_i is assumed to be a parametric probabilistic model: let $p_Y(y_i; \psi_i)$ be the probability distribution of y_i , where ψ_i is a vector of parameters.

In a population framework, the vector of parameters ψ_i is assumed to be drawn from a population distribution $p_\Psi(\psi_i; \theta)$ where θ is a vector of population parameters.

Then, the probabilistic model is the joint probability distribution

$$p(y_i, \psi_i; \theta) = p_Y(y_i | \psi_i) p_\Psi(\psi_i; \theta) \quad (1)$$

To define a model thus consists in defining precisely these two terms.

In most applications, the observed data y_i are continuous longitudinal data. We then assume the following representation for y_i :

$$y_{ij} = f(t_{ij}, \psi_i) + g(t_{ij}, \psi_i) \varepsilon_{ij} \quad , \quad 1 \leq i \leq N \quad , \quad 1 \leq j \leq n_i. \quad (2)$$

Here, y_{ij} is the observation obtained from subject i at time t_{ij} . The residual errors (ε_{ij}) are assumed to be standardized random variables (mean zero and variance 1). The residual error model is represented by function g in model (2).

Function f is usually the solution to a system of ordinary differential equations (pharmacokinetic/pharmacodynamic models, etc.) or a system of partial differential equations (tumor growth, respiratory system, etc.). This component is a fundamental component of the model since it defines the prediction of the observed kinetics for a given set of parameters.

The vector of individual parameters ψ_i is usually function of a vector of population parameters ψ_{pop} , a vector of random effects $\eta_i \sim \mathcal{N}(0, \Omega)$, a vector of individual covariates c_i (weight, age, gender, ...) and some fixed effects β .

The joint model of y and ψ depends then on a vector of parameters $\theta = (\psi_{\text{pop}}, \beta, \Omega)$.

3.3. Computational Statistical Methods

Central to modern statistics is the use of probabilistic models. To relate these models to data requires the ability to calculate the probability of the observed data: the likelihood function, which is central to most statistical methods and provides a principled framework to handle uncertainty.

The emergence of computational statistics as a collection of powerful and general methodologies for carrying out likelihood-based inference made complex models with non-standard data accessible to likelihood, including hierarchical models, models with intricate latent structure, and missing data.

In particular, algorithms previously developed by POPIX for mixed effects models, and today implemented in several software tools (especially MONOLIX) are part of these methods:

- the adaptive Metropolis-Hastings algorithm allows one to sample from the conditional distribution of the individual parameters $p(\psi_i | y_i; c_i, \theta)$,
- the SAEM algorithm is used to maximize the observed likelihood $\mathcal{L}(\theta; y) = p(y; \theta)$,
- Importance Sampling Monte Carlo simulations provide an accurate estimation of the observed log-likelihood $\log(\mathcal{L}(\theta; y))$.

Computational statistics is an area which remains extremely active today. Recently, one can notice that the incentive for further improvements and innovation comes mainly from three broad directions: the high dimensional challenge, the quest for adaptive procedures that can eliminate the cumbersome process of tuning "by hand" the settings of the algorithms and the need for flexible theoretical support, arguably required by all recent developments as well as many of the traditional MCMC algorithms that are widely used in practice.

Working in these three directions is a clear objective for XPOP.

3.4. Markov Chain Monte Carlo algorithms

While these Monte Carlo algorithms have turned into standard tools over the past decade, they still face difficulties in handling less regular problems such as those involved in deriving inference for high-dimensional models. One of the main problems encountered when using MCMC in this challenging settings is that it is difficult to design a Markov chain that efficiently samples the state space of interest.

The Metropolis-adjusted Langevin algorithm (MALA) is a Markov chain Monte Carlo (MCMC) method for obtaining random samples from a probability distribution for which direct sampling is difficult. As the name suggests, MALA uses a combination of two mechanisms to generate the states of a random walk that has the target probability distribution as an invariant measure:

1. new states are proposed using Langevin dynamics, which use evaluations of the gradient of the target probability density function;
2. these proposals are accepted or rejected using the Metropolis-Hastings algorithm, which uses evaluations of the target probability density (but not its gradient).

Informally, the Langevin dynamics drives the random walk towards regions of high probability in the manner of a gradient flow, while the Metropolis-Hastings accept/reject mechanism improves the mixing and convergence properties of this random walk.

Several extensions of MALA have been proposed recently by several authors, including fMALA (fast MALA), AMALA (anisotropic MALA), MMALA (manifold MALA), position-dependent MALA (PMALA), ...

MALA and these extensions have demonstrated to represent very efficient alternative for sampling from high dimensional distributions. We therefore need to adapt these methods to general mixed effects models.

3.5. Parameter estimation

The Stochastic Approximation Expectation Maximization (SAEM) algorithm has shown to be extremely efficient for maximum likelihood estimation in incomplete data models, and particularly in mixed effects models for estimating the population parameters. However, there are several practical situations for which extensions of SAEM are still needed:

High dimensional model: a complex physiological model may have a large number of parameters (in the order of 100). Then several problems arise:

- when most of these parameters are associated with random effects, the MCMC algorithm should be able to sample, for each of the N individuals, parameters from a high dimensional distribution. Efficient MCMC methods for high dimensions are then required.
- Practical identifiability of the model is not ensured with a limited amount of data. In other words, we cannot expect to be able to properly estimate all the parameters of the model, including the fixed effects and the variance-covariance matrix of the random effects. Then, some random effects should be removed, assuming that some parameters do not vary in the population. It may also be necessary to fix the value of some parameters (using values from the literature for instance). The strategy to decide which parameters should be fixed and which random effects should be removed remains totally empirical. XPOP aims to develop a procedure that will help the modeller to take such decisions.

Large number of covariates: the covariate model aims to explain part of the inter-patient variability of some parameters. Classical methods for covariate model building are based on comparisons with respect to some criteria, usually derived from the likelihood (AIC, BIC), or some statistical test (Wald test, LRT, etc.). In other words, the modelling procedure requires two steps: first, all possible models are fitted using some estimation procedure (e.g. the SAEM algorithm) and the likelihood of each model is computed using a numerical integration procedure (e.g. Monte Carlo Importance Sampling); then, a model selection procedure chooses the "best" covariate model. Such a strategy is only possible with a reduced number of covariates, i.e., with a "small" number of models to fit and compare.

As an alternative, we are thinking about a Bayesian approach which consists of estimating simultaneously the covariate model and the parameters of the model in a single run. An (informative or uninformative) prior is defined for each model by defining a prior probability for each covariate to be included in the model. In other words, we extend the probabilistic model by introducing binary variables that indicate the presence or absence of each covariate in the model. Then, the model selection procedure consists of estimating and maximizing the conditional distribution of this sequence of binary variables. Furthermore, a probability can be associated to any of the possible covariate models.

This conditional distribution can be estimated using an MCMC procedure combined with the SAEM algorithm for estimating the population parameters of the model. In practice, such an approach can only deal with a limited number of covariates since the dimension of the probability space to explore increases exponentially with the number of covariates. Consequently, we would like to have methods able to find a small number of variables (from a large starting set) that influence certain parameters in populations of individuals. That means that, instead of estimating the conditional distribution of all the covariate models as described above, the algorithm should focus on the most likely ones.

Fixed parameters: it is quite frequent that some individual parameters of the model have no random component and are purely fixed effects. Then, the model may not belong to the exponential family anymore and the original version of SAEM cannot be used as it is. Several extensions exist:

- introduce random effects with decreasing variances for these parameters,
- introduce a prior distribution for these fixed effects,
- apply the stochastic approximation directly on the sequence of estimated parameters, instead of the sufficient statistics of the model.

None of these methods always work correctly. Furthermore, what are the pros and cons of these methods is not clear at all. Then, developing a robust methodology for such model is necessary.

Convergence toward the global maximum of the likelihood: convergence of SAEM can strongly depend on the initial guess when the observed likelihood has several local maxima. A kind of simulated annealing version of SAEM was previously developed and implemented in MONOLIX. The method works quite well in most situations but there is no theoretical justification and choosing the settings of this algorithm (i.e. how the temperature decreases during the iterations) remains empirical. A precise analysis of the algorithm could be very useful to better understand why it "works" in practice and how to optimize it.

Convergence diagnostic: Convergence of SAEM was theoretically demonstrated under very general hypothesis. Such result is important but of little interest in practice at the time to use SAEM in a finite amount of time, i.e. in a finite number of iterations. Some qualitative and quantitative criteria should be defined in order to both optimize the settings of the algorithm, detect a poor convergence of SAEM and evaluate the quality of the results in order to avoid using them unwisely.

3.6. Model building

Defining an optimal strategy for model building is far from easy because a model is the assembled product of numerous components that need to be evaluated and perhaps improved: the structural model, residual error model, covariate model, covariance model, etc.

How to proceed so as to obtain the best possible combination of these components? There is no magic recipe but an effort will be made to provide some qualitative and quantitative criteria in order to help the modeller for building his model.

The strategy to take will mainly depend on the time we can dedicate to building the model and the time required for running it. For relatively simple models for which parameter estimation is fast, it is possible to fit many models and compare them. This can also be done if we have powerful computing facilities available (e.g., a cluster) allowing large numbers of simultaneous runs.

However, if we are working on a standard laptop or desktop computer, model building is a sequential process in which a new model is tested at each step. If the model is complex and requires significant computation time (e.g., when involving systems of ODEs), we are constrained to limit the number of models we can test in a reasonable time period. In this context, it also becomes important to carefully choose the tasks to run at each step.

3.7. Model evaluation

Diagnostic tools are recognized as an essential method for model assessment in the process of model building. Indeed, the modeler needs to confront "his" model with the experimental data before concluding that this model is able to reproduce the data and before using it for any purpose, such as prediction or simulation for instance.

The objective of a diagnostic tool is twofold: first we want to check if the assumptions made on the model are valid or not ; then, if some assumptions are rejected, we want to get some guidance on how to improve the model.

As is the usual case in statistics, it is not because this "final" model has not been rejected that it is necessarily the "true" one. All that we can say is that the experimental data does not allow us to reject it. It is merely one of perhaps many models that cannot be rejected.

Model diagnostic tools are for the most part graphical, i.e., visual; we "see" when something is not right between a chosen model and the data it is hypothesized to describe. These diagnostic plots are usually based on the empirical Bayes estimates (EBEs) of the individual parameters and EBEs of the random effects: scatterplots of individual parameters versus covariates to detect some possible relationship, scatterplots of pairs of random effects to detect some possible correlation between random effects, plot of the empirical distribution of the random effects (boxplot, histogram,...) to check if they are normally distributed, ...

The use of EBEs for diagnostic plots and statistical tests is efficient with rich data, i.e. when a significant amount of information is available in the data for recovering accurately all the individual parameters. On the contrary, tests and plots can be misleading when the estimates of the individual parameters are greatly shrunk.

We propose to develop new approaches for diagnosing mixed effects models in a general context and derive formal and unbiased statistical tests for testing separately each feature of the model.

3.8. Missing data

The ability to easily collect and gather a large amount of data from different sources can be seen as an opportunity to better understand many processes. It has already led to breakthroughs in several application areas. However, due to the wide heterogeneity of measurements and objectives, these large databases often exhibit an extraordinary high number of missing values. Hence, in addition to scientific questions, such data also present some important methodological and technical challenges for data analyst.

Missing values occur for a variety of reasons: machines that fail, survey participants who do not answer certain questions, destroyed or lost data, dead animals, damaged plants, etc. Missing values are problematic since most statistical methods can not be applied directly on a incomplete data. Many progress have been made to properly handle missing values. However, there are still many challenges that need to be addressed in the future, that are crucial for the users.

- State of arts methods often consider the case of continuous or categorical data whereas real data are very often mixed. The idea is to develop a multiple imputation method based on a specific principal component analysis (PCA) for mixed data. Indeed, PCA has been used with success to predict (impute) the missing values. A very appealing property is the ability of the method to handle very large matrices with large amount of missing entries.
- The asymptotic regime underlying modern data is not any more to consider that the sample size increases but that both number of observations and number of variables are very large. In practice first experiments showed that the coverage properties of confidence areas based on the classical methods to estimate variance with missing values varied widely. The asymptotic method and the bootstrap do well in low-noise setting, but can fail when the noise level gets high or when the number of variables is much greater than the number of rows. On the other hand, the jackknife has good coverage properties for large noisy examples but requires a minimum number of variables to be stable enough.
- Inference with missing values is usually performed under the assumption of "Missing at Random" (MAR) values which means that the probability that a value is missing may depend on the observed data but does not depend on the missing value itself. In real data and in particular in data coming from clinical studies, both "Missing Non at Random" (MNAR) and MAR values occur. Taking into account in a proper way both types of missing values is extremely challenging but is worth investigating since the applications are extremely broad.

It is important to stress that missing data models are part of the general incomplete data models addressed by XPOP. Indeed, models with latent variables (i.e. non observed variables such as random effects in a mixed effects model), models with censored data (e.g. data below some limit of quantification) or models with dropout mechanism (e.g. when a subject in a clinical trial fails to continue in the study) can be seen as missing data models.

4. Application Domains

4.1. Precision medicine and pharmacogenomics

Pharmacogenomics involves using an individual's genome to determine whether or not a particular therapy, or dose of therapy, will be effective. Indeed, people's reaction to a given drug depends on their physiological state and environmental factors, but also to their individual genetic make-up.

Precision medicine is an emerging approach for disease treatment and prevention that takes into account individual variability in genes, environment, and lifestyle for each person. While some advances in precision medicine have been made, the practice is not currently in use for most diseases.

Currently, in the traditional population approach, inter-individual variability in the reaction to drugs is modeled using covariates such as weight, age, sex, ethnic origin, etc. Genetic polymorphisms susceptible to modify pharmacokinetic or pharmacodynamic parameters are much harder to include, especially as there are millions of possible polymorphisms (and thus covariates) per patient.

The challenge is to determine which genetic covariates are associated to some PKPD parameters and/or implicated in patient responses to a given drug.

Another problem encountered is the dependence of genes, as indeed, gene expression is a highly regulated process. In cases where the explanatory variables (genomic variants) are correlated, Lasso-type methods for model selection are thwarted.

There is therefore a clear need for new methods and algorithms for the estimation, validation and selection of mixed effects models adapted to the problems of genomic medicine.

A target application of this project concerns the lung cancer.

EGFR (Epidermal Growth Factor Receptor) is a cell surface protein that binds to epidermal growth factor. We know that deregulation of the downstream signaling pathway of EGFR is involved in the development of lung cancers and several gene mutations responsible for this deregulation are known.

Our objective is to identify the variants responsible for the disruption of this pathway using a modelling approach. The data that should be available for developing such model are ERK (Extracellular signal-regulated kinases) phosphorylation time series, obtained from different genetic profiles.

The model that we aim to develop will describe the relationship between the parameters of the pathway and the genomic covariates, i.e. the genetic profile. Variants related to the pathway include: variants that modify the affinity binding of ligands to receptors, variants that modify the total amount of protein, variants that affect the catalytic site,...

4.2. Oncology

In cancer, the most dreadful event is the formation of metastases that disseminate tumor cells throughout the organism. Cutaneous melanoma is a cancer, where the primary tumor can easily be removed by surgery. However, this cancer is of poor prognosis; because melanomas metastasize often and rapidly. Many melanomas arise from excessive exposure to mutagenic UV from the sun or sunbeds. As a consequence, the mutational burden of melanomas is generally high

RAC1 encodes a small GTPase that induces cell cycle progression and migration of melanoblasts during embryonic development. Patients with the recurrent P29S mutation of RAC1 have 3-fold increased odds at having regional lymph nodes invaded at the time of diagnosis. RAC1 is unlikely to be a good therapeutic target, since a potential inhibitor that would block its catalytic activity, would also lock it into the active GTP-bound state. This project thus investigates the possibility of targeting the signaling pathway downstream of RAC1.

XPOP is mainly involved in Task 1 of the project: *Identifying deregulations and mutations of the ARP2/3 pathway in melanoma patients.*

Association of over-expression or down-regulation of each marker with poor prognosis in terms of invasion of regional lymph nodes, metastases and survival, will be examined using classical univariate and multivariate analysis. We will then develop specific statistical models for survival analysis in order to associate prognosis factors to each composition of complexes. Indeed, one has to implement the further constraint that each subunit has to be contributed by one of several paralogous subunits. An original method previously developed by XPOP has already been successfully applied to WAVE complex data in breast cancer.

The developed models will be rendered user-friendly through a dedicated Rsoftware package.

This project can represent a significant step forward in precision medicine of the cutaneous melanoma.

4.3. Hemodialysis

Hemodialysis is a process for removing waste and excess water from the blood and is used primarily as an artificial replacement for lost kidney function in people with kidney failure. Side effects caused by removing too much fluid and/or removing fluid too rapidly include low blood pressure, fatigue, chest pains, leg-cramps, nausea and headaches.

Nephrologists must therefore correctly assess the hydration status in chronic hemodialysis patients and consider fluid overload effects when prescribing dialysis, according to a new study.

The fluid overload biomarker, B-type natriuretic peptide (BNP) is an important component of managing patients with kidney disease. Indeed, it is believed that each dialysis patient will have an ideal or "dry" BNP level which will accurately and reproducibly reflect their optimal fluid status.

The objective of this study is to develop a model for the BNP and the hydration status using individual information (age, sex, ethnicity, systolic blood pressure, BMI, coronary heart disease history, ...).

The impact will be significant if the method succeeds. Indeed, it will be possible for the nephrologists to use this model for monitoring individually each treatment, in order to avoid risks of hypotension (low BNP) or overweight (high BNP).

4.4. Intracellular processes

Significant cell-to-cell heterogeneity is ubiquitously-observed in isogenic cell populations. Cells respond differently to a same stimulation. For example, accounting for such heterogeneity is essential to quantitatively understand why some bacteria survive antibiotic treatments, some cancer cells escape drug-induced suicide, stem cell do not differentiate, or some cells are not infected by pathogens.

The origins of the variability of biological processes and phenotypes are multifarious. Indeed, the observed heterogeneity of cell responses to a common stimulus can originate from differences in cell phenotypes (age, cell size, ribosome and transcription factor concentrations, etc), from spatio-temporal variations of the cell environments and from the intrinsic randomness of biochemical reactions. From systems and synthetic biology perspectives, understanding the exact contributions of these different sources of heterogeneity on the variability of cell responses is a central question.

The main ambition of this project is to propose a paradigm change in the quantitative modelling of cellular processes by shifting from mean-cell models to single-cell and population models. The main contribution of XPOP focuses on methodological developments for mixed-effects model identification in the context of growing cell populations.

- Mixed-effects models usually consider an homogeneous population of independent individuals. This assumption does not hold when the population of cells (i.e. the statistical individuals) consists of several generations of dividing cells. We then need to account for inheritance of single-cell parameters in this population. More precisely, the problem is to attribute the new state and parameter values to newborn cells given (the current estimated values for) the mother.
- The mixed-effects modelling framework corresponds to a strong assumption: differences between cells are static in time (ie, cell-specific parameters have fixed values). However, it is likely that for any given cell, ribosome levels slowly vary across time, since like any other protein, ribosomes are produced in a stochastic manner. We will therefore extend our modelling framework so as to account for the possible random fluctuations of parameter values in individual cells. Extensions based on stochastic differential equations will be investigated.
- Identifiability is a fundamental prerequisite for model identification and is also closely connected to optimal experimental design. We will derive criteria for theoretical identifiability, in which different parameter values lead to non-identical probability distributions, and for structural identifiability, which concerns the algebraic properties of the structural model, i.e. the ODE system. We will then address the problem of practical identifiability, whereby the model may be theoretically identifiable but the design of the experiment may make parameter estimation difficult and imprecise. An interesting problem is whether accounting for lineage effects can help practical identifiability of the parameters of the individuals in presence of measurement and biological noise.

4.5. Population pharmacometrics

Pharmacometrics involves the analysis and interpretation of data produced in pre-clinical and clinical trials. Population pharmacokinetics studies the variability in drug exposure for clinically safe and effective doses by focusing on identification of patient characteristics which significantly affect or are highly correlated with this variability. Disease progress modeling uses mathematical models to describe, explain, investigate and predict the changes in disease status as a function of time. A disease progress model incorporates functions describing natural disease progression and drug action.

The model based drug development (MBDD) approach establishes quantitative targets for each development step and optimizes the design of each study to meet the target. Optimizing study design requires simulations, which in turn require models. In order to arrive at a meaningful design, mechanisms need to be understood and correctly represented in the mathematical model. Furthermore, the model has to be predictive for future studies. This requirement precludes all purely empirical modeling; instead, models have to be mechanistic.

In particular, physiologically based pharmacokinetic models attempt to mathematically transcribe anatomical, physiological, physical, and chemical descriptions of phenomena involved in the ADME (Absorption - Distribution - Metabolism - Elimination) processes. A system of ordinary differential equations for the quantity of substance in each compartment involves parameters representing blood flow, pulmonary ventilation rate, organ volume, etc.

The ability to describe variability in pharmacometrics model is essential. The nonlinear mixed-effects modeling approach does this by combining the structural model component (the ODE system) with a statistical model, describing the distribution of the parameters between subjects and within subjects, as well as quantifying the unexplained or residual variability within subjects.

The objective of XPOP is to develop new methods for models defined by a very large ODE system, a large number of parameters and a large number of covariates. Contributions of XPOP in this domain are mainly methodological and there is no privileged therapeutic application at this stage.

However, it is expected that these new methods will be implemented in software tools, including MONOLIX and Rpackages for practical use.

4.6. Mass spectrometry

One of the main recent developments in analytical chemistry is the rapid democratization of high-resolution mass spectrometers. These instruments produce extremely complex mass spectra, which can include several hundred thousand ions when analyzing complex samples. The analysis of complex matrices (biological, agri-food, cosmetic, pharmaceutical, environmental, etc.) is precisely one of the major analytical challenges of this new century. Academic and industrial researchers are particularly interested in trying to quickly and effectively establish the chemical consequences of an event on a complex matrix. This may include, for example, searching for pesticide degradation products and metabolites in fruits and vegetables, photoproducts of active ingredients in a cosmetic emulsion exposed to UV rays or chlorination products of biocides in hospital effluents. The main difficulty of this type of analysis is based on the high spatial and temporal variability of the samples, which is in addition to the experimental uncertainties inherent in any measurement and requires a large number of samples and analyses to be carried out and computerized data processing (up to 16 million per mass spectrum).

A collaboration between XPOP and the Molecular Chemistry Laboratory (LCM) of the Ecole Polytechnique began in 2018. Our objective is to develop new methods for the statistical analysis of mass spectrometry data.

5. Highlights of the Year

5.1. Highlights of the Year

Version 1.0 of the SPIX software was available in November 2018.

6. New Software and Platforms

6.1. mlxR

KEYWORDS: Simulation - Data visualization - Clinical trial simulator

FUNCTIONAL DESCRIPTION: The models are encoded using the model coding language 'Mlxtran', automatically converted into C++ codes, compiled on the fly and linked to R using the 'Rcpp' package. That allows one to implement very easily complex ODE-based models and complex statistical models, including mixed effects models, for continuous, count, categorical, and time-to-event data.

- Contact: Marc Lavielle
- URL: <http://simulx.webpopix.org/>

6.2. Rsmlx

R speaks Monolix

KEYWORDS: Data modeling - Nonlinear mixed effects models - Statistical modeling

FUNCTIONAL DESCRIPTION: Among other tasks, 'Rsmlx' provides a powerful tool for automatic PK model building, performs statistical tests for model assessment, bootstrap simulation and likelihood profiling for computing confidence intervals. 'Rsmlx' also proposes several automatic covariate search methods for mixed effects models.

- Partner: Lixoft
- Contact: Marc Lavielle

6.3. SPIX

KEYWORDS: Data modeling - Mass spectrometry - Chemistry

FUNCTIONAL DESCRIPTION: SPIX allows you to - To automatically identify, on the basis of statistical approaches, small but significant differences in spectra measured under different conditions, - Model the kinetics of entities that evolve over time

- Partner: Laboratoire de Chimie Moléculaire - Ecole Polytechnique
- Contact: Marc Lavielle

7. New Results

7.1. Normalizing constants of log-concave densities

We derive explicit bounds for the computation of normalizing constants Z for log-concave densities $\pi = e^{-U}/Z$ w.r.t. the Lebesgue measure on \mathbb{R}^d . Our approach relies on a Gaussian annealing combined with recent and precise bounds on the Unadjusted Langevin Algorithm (Durmus, A. and Moulines, E. (2016). High-dimensional Bayesian inference via the Unadjusted Langevin Algorithm). Polynomial bounds in the dimension d are obtained with an exponent that depends on the assumptions made on U . The algorithm also provides a theoretically grounded choice of the annealing sequence of variances. A numerical experiment supports our findings. Results of independent interest on the mean squared error of the empirical average of locally Lipschitz functions are established.

7.2. The Tamed Unadjusted Langevin Algorithm

We consider the problem of sampling from a probability measure π having a density on \mathbb{R}^d known up to a normalizing constant, $x \rightarrow e^{-U(x)}/Z$. The Euler discretization of the Langevin stochastic differential equation (SDE) is known to be unstable in a precise sense, when the potential U is superlinear. Based on previous works on the taming of superlinear drift coefficients for SDEs, we introduce the Tamed Unadjusted Langevin Algorithm (TULA) and obtain non-asymptotic bounds in V -total variation norm and Wasserstein distance of order 2 between the iterates of TULA and π , as well as weak error bounds. Numerical experiments support our findings.

7.3. Development and performance of npde for the evaluation of time-to-event models

Normalised prediction distribution errors (npde) are used to graphically and statistically evaluate mixed-effect models for continuous responses. Our aim was to extend npde to time-to-event (TTE) models and evaluate their performance. We extended npde to TTE models using imputations to take into account censoring. We then evaluated their performance in terms of type I error and power to detect model misspecifications for TTE data by means of a simulation study with different sample sizes. Type I error was found to be close to the expected 5% significance level for all sample sizes tested. The npde were able to detect misspecifications in the baseline hazard as well as in the link between the longitudinal variable and the survival function. The ability to detect model misspecifications increased as the difference in the shape of the survival function became more apparent. As expected, the power also increased as the sample size increased. Imputing the censored events tended to decrease the percentage of rejections.

7.4. Low-rank Interaction with Sparse Additive Effects Model for Large Data Frames

Many applications of machine learning involve the analysis of large data frames-matrices collecting heterogeneous measurements (binary, numerical, counts, etc.) across samples-with missing values. Low-rank models are popular in this framework for tasks such as visualization, clustering and missing value imputation. Yet, available methods with statistical guarantees and efficient optimization do not allow explicit modeling of main additive effects such as row and column, or covariate effects. We introduced a low-rank interaction and sparse additive effects (LORIS) model which combines matrix regression on a dictionary and low-rank design, to estimate main effects and interactions simultaneously. We provide statistical guarantees in the form of upper bounds on the estimation error of both components. Then, we introduced a mixed coordinate gradient descent (MCGD) method which provably converges sub-linearly to an optimal solution and is computationally efficient for large scale data sets. Simulated and survey data showed that the method has a clear advantage over current practices, which consist in dealing separately with additive effects in a preprocessing step.

7.5. Diffusion approximations and control variates for MCMC

A new methodology was developed for the construction of control variates to reduce the variance of additive functionals of Markov Chain Monte Carlo (MCMC) samplers. Our control variates are defined as linear combinations of functions whose coefficients are obtained by minimizing a proxy for the asymptotic variance. The construction is theoretically justified by two new results. We first show that the asymptotic variances of some well-known MCMC algorithms, including the Random Walk Metropolis and the (Metropolis) Unadjusted/Adjusted Langevin Algorithm, are close to the asymptotic variance of the Langevin diffusion. Second, we provide an explicit representation of the optimal coefficients minimizing the asymptotic variance of the Langevin diffusion. Several examples of Bayesian inference problems demonstrate that the corresponding reduction in the variance is significant, and that in some cases it can be dramatic.

7.6. Density estimation for random walks in random environment

We consider the problem of non-parametric density estimation of a random environment from the observation of a single trajectory of a random walk in this environment. We first construct a density estimator using the beta-moments. We then show that the Goldenshluger-Lepski method can be used to select the beta-moment. We prove nonasymptotic bounds for the supremum norm of these estimators for both the recurrent and the transient to the right case. A simulation study supports our theoretical findings.

7.7. Imputation of mixed data with multilevel singular value decomposition

Statistical analysis of large data sets offers new opportunities to better understand many processes. Yet, data accumulation often implies relaxing acquisition procedures or compounding diverse sources. As a consequence, such data sets often contain mixed data, i.e. both quantitative and qualitative and many missing values. Furthermore, aggregated data present a natural multilevel structure, where individuals or samples are nested within different sites, such as countries or hospitals. Imputation of multilevel data has therefore drawn some attention recently, but current solutions are not designed to handle mixed data, and suffer from important drawbacks such as their computational cost. In this article, we propose a single imputation method for multilevel data, which can be used to complete either quantitative, categorical or mixed data. The method is based on multilevel singular value decomposition (SVD), which consists in decomposing the variability of the data into two components, the between and within groups variability, and performing SVD on both parts. We show on a simulation study that in comparison to competitors, the method has the great advantages of handling data sets of various size, and being computationally faster. Furthermore, it is the first so far to handle mixed data. We apply the method to impute a medical data set resulting from the aggregation of several data sets coming from different hospitals. This application falls in the framework of a larger project on Trauma patients. To overcome obstacles associated to the aggregation of medical data, we turn to distributed computation. The method is implemented in an R package

7.8. Logistic Regression with Missing Covariates – Parameter Estimation, Model Selection and Prediction

Logistic regression is a common classification method in supervised learning. Surprisingly, there are very few solutions for performing it and selecting variables in the presence of missing values. We develop a complete approach, including the estimation of parameters and variance of estimators, derivation of confidence intervals and a model selection procedure, for cases where the missing values can be anywhere in covariates. By well organizing different patterns of missingness in each observation, we propose a stochastic approximation version of the EM algorithm based on Metropolis-Hastings sampling, to perform statistical inference for logistic regression with incomplete data. We also tackle the problem of prediction for a new individual with missing values, which is never addressed. The methodology is computationally efficient, and its good coverage and variable selection properties are demonstrated in a simulation study where we contrast its performances to other methods. For instance, the popular multiple imputation by chained equation can lead to biased estimates while our method is unbiased. The method was applied on a dataset of severely traumatized patients from Paris hospitals to predict the occurrence of hemorrhagic shock, a leading cause of early preventable death in severe trauma cases. The aim is to consolidate the current red flag procedure, a binary alert identifying patients with a high risk of severe hemorrhage. The methodology is implemented in the R package *misaem*.

7.9. A fast Stochastic Approximation of the EM algorithm for nonlinear mixed effects models

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 demonstrate the good 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.

7.10. Incomplete graphical model inference via latent tree aggregation

Graphical network inference is used in many fields such as genomics or ecology to infer the conditional independence structure between variables, from measurements of gene expression or species abundances for instance. In many practical cases, not all variables involved in the network have been observed, and the samples are actually drawn from a distribution where some variables have been marginalized out. This challenges the sparsity assumption commonly made in graphical model inference, since marginalization yields locally dense structures, even when the original network is sparse. We developed a procedure for inferring Gaussian graphical models when some variables are unobserved, that accounts both for the influence of missing variables and the low density of the original network. Our model is based on the aggregation of spanning trees, and the estimation procedure on the Expectation-Maximization algorithm. We treat the graph structure and the unobserved nodes as missing variables and compute posterior probabilities of edge appearance. To provide a complete methodology, we also propose several model selection criteria to estimate the number of missing nodes. A simulation study and an illustration flow cytometry data reveal that our method has favorable edge detection properties compared to existing graph inference techniques. The methods are implemented in an R package.

8. Bilateral Contracts and Grants with Industry

8.1. Bilateral Contracts with Industry

Contract with Dassault Systèmes

9. Partnerships and Cooperations

9.1. National Initiatives

9.1.1. ANR

Mixed-Effects Models of Intracellular Processes: Methods, Tools and Applications (MEMIP)

Coordinator: Gregory Batt (InBio Inria team)

Other partners: InBio and IBIS Inria teams, Laboratoire Matière et Systèmes Complexes (UMR 7057; CNRS and Paris Diderot Univ.)

9.1.2. Institut National du Cancer (INCa)

Targeting Rac-dependent actin polymerization in cutaneous melanoma - Institut National du Cancer

Coordinator: Alexis Gautreau (Ecole Polytechnique)

Other partners: Laboratoire de Biochimie (Polytechnique), Institut Curie, INSERM.

9.2. International Initiatives

9.2.1. International Initiatives

SaSMoTiDep

Title: Statistical and Stochastic modeling for time-dependent data

International Partners (Institution - Laboratory - Researcher):

Universidad de Valparaíso (Chile) - Centro de Investigación y Modelamiento de Fenómenos Aleatorios Valparaíso (CIMFAV) - Cristian Meza Becerra

Universidad Nacional de Colombia (Colombia) - Department of Statistics - Viswanathan Arunachalam

Duration: 01/01/2018 - 31/12/2019

Start year: 2018

See also: <https://sasmotidep.uv.cl>

In many applications, multiple measurements are made on one or several experimental units over a period of time. Such data could be called time-dependent data. From a statistical point of view, if we consider only one experimental unit, we can use a time series analysis. In the other hand, if we consider experimental designs (or observational studies) for several experimental units (or subjects) where each subject is measured at several points in time, we can use the term longitudinal data. In this project, we propose to study several statistical and stochastic models for repeated measures using parametric and non-parametric approaches. In particular, we will study the inference in complex mixed effects models, we will propose novel segmentation models for multiple series, non-parametric methods in dependent models and stochastic models. We will apply these methods to real data from several fields as biometrics, reliability, population dynamics and finance.

9.3. International Research Visitors

9.3.1. Visits of International Scientists

Ricardo Rios, Universidad Central de Venezuela, Caracas: September 2018.

10. Dissemination

10.1. Promoting Scientific Activities

10.1.1. Scientific Events Organisation

10.1.1.1. General Chair, Scientific Chair

- Winter School [New Frontiers in high-dimensional probability and statistics](#), Moscow, February 23-24, 2018
- Workshop "[Structural Inference in Inference in High-Dimensional Models](#)", Moscow, 5-8 September 2018

10.1.2. Scientific Events Selection

10.1.2.1. Member of the Conference Program Committees

- [NIPS 2018](#)
- [COLT 2018](#)
- [useR!2018](#)

10.1.3. Journal

10.1.3.1. Member of the Editorial Boards

- Stochastic Processes and their Applications
- Journal of Statistical Planning and Inference
- Journal of Computational and Graphical Statistics

10.1.4. Invited Talks

- [Analysis of adaptive stochastic gradient and MCMC algorithms](#), June 2018
- [Langevin MCMC: theory and methods](#), Singapour, June 2018

- **MCQMC 2018: Langevin Monte Carlo: theory and methods**, Rennes, July 2018
- **Langevin MCMC: theory and methods**, Singapour, August 2018
- **Huawei European Innovation Day, "Toward Bayesian Machine Learning"**, Roma, November 2018
- **Iberoamerican Pharmacometrics Network Congress 2018**, Guadalajara, Mexico.
- "Tests d'hypothèse et algorithmes de construction de modèle pour les modèles à effets mixtes", Toulouse, September 2018.
- **New algorithms and tools for population model building and evaluation**, Valparaiso, Chile, December 2018.

10.1.5. Leadership within the Scientific Community

- Eric Moulines is in charge of the academic supervision of the **International Laboratory of Stochastic Algorithms and High-dimensional inference**, National Research University, Higher School of Economics, funded by the Russian Academic Excellence Project.
- Eric Moulines is associate researcher of the **Alan Turing Institute**
- Eric Moulines is elected member of the French Académie des Sciences.

10.1.6. Scientific Expertise

- Marc Lavielle is member of the Scientific Committee of the High Council for Biotechnologies.
- Marc Lavielle is member of the evaluation committee of the **International Center for Mathematics (CIMAT)**, Guanajuato, Mexico.
- Eric Moulines is member of the award committee of foundation "Charles Defforey".

10.1.7. Research administration

- Marc Lavielle is member of the Scientific Programming Committee (CPS) of the Institute Henri Poincaré (IHP).
- Eric Moulines is a board member of the Institut de Convergence DataIA.
- Julie Josse is elected member of the R foundation.

10.2. Teaching - Supervision - Juries

10.2.1. Teaching

Master : Julie Josse, Statistics with R, 48, M2, X-HEC
 Master : Eric Moulines, Regression models, 36, M2, X-HEC
 Engineering School : Eric Moulines, Statistics, 36, 2A, X
 Engineering School : Eric Moulines, Markov Chains, 36, 3A, X
 Engineering School : Erwan Le Pennec, Statistics, 36, 2A, X
 Engineering School : Erwan Le Pennec, Statistical Learning, 36, 3A, X
 Engineering School : Marc Lavielle, Statistics in Action, 48, 3A, X

10.2.2. Supervision

PhD in progress : Nicolas Brosse, September 2016, Eric Moulines
 PhD in progress : Geneviève Robin, September 2016, Julie Josse and Eric Moulines
 PhD in progress : Belhal Karimi, October 2016, Marc Lavielle and Eric Moulines
 PhD in progress : Marine Zulian, October 2016, Marc Lavielle
 PhD in progress : Wei Jiang , October 2017, Julie Josse and Marc Lavielle

10.3. Popularization

10.3.1. Internal or external Inria responsibilities

Eric Moulines is member of the Evaluation Committee of Inria.

10.3.2. Creation of media or tools for science outreach

Marc Lavielle developed and maintains the learning platform [Statistics in Action](#). The purpose of this online learning platform is to show how statistics (and biostatistics) may be efficiently used in practice using R. It is specifically geared towards teaching statistical modelling concepts and applications for self-study. Indeed, most of the available teaching material tends to be quite "static" while statistical modelling is very much subject to "learning by doing".

Julie Josse (with Nicholas Tierney and Nathalie Vialaneix) developed and maintains the website [R-miss-tastic](#), A resource website on missing data.

11. Bibliography

Publications of the year

Articles in International Peer-Reviewed Journals

- [1] F. ALLA, N. AGRINIER, M. LAVIELLE, P. ROSSIGNOL, D. GONTHIER, J.-M. BOIVIN, F. ZANNAD. *Impact of the interruption of a large heart failure regional disease management program on hospital admission rate: a population-based study*, in "European Journal of Heart Failure", 2018, <https://hal.inria.fr/hal-01727632>
- [2] N. BROSSE, A. DURMUS, É. MOULINES. *Normalizing constants of log-concave densities*, in "Electronic journal of statistics ", March 2018, vol. 12, n^o 1, pp. 851-889 [DOI : 10.1214/18-EJS1411], <https://hal.inria.fr/hal-01648666>
- [3] N. BROSSE, A. DURMUS, É. MOULINES, S. SABANIS. *The Tamed Unadjusted Langevin Algorithm*, in "Stochastic Processes and their Applications", October 2018 [DOI : 10.1016/J.SPA.2018.10.002], <https://hal.inria.fr/hal-01648667>
- [4] M. CEROU, M. LAVIELLE, K. BRENDEL, M. CHENEL, E. COMETS. *Development and performance of npde for the evaluation of time-to-event models*, in "Pharmaceutical Research", January 2018, vol. 35, n^o 2, 30 p. [DOI : 10.1007/s11095-017-2291-3], <https://www.hal.inserm.fr/inserm-01695500>
- [5] M. LAVIELLE. *Pharmacometrics Models with Hidden Markovian Dynamics*, in "Journal of Pharmacokinetics and Pharmacodynamics", 2018, vol. 45, n^o 1, pp. 91–105 [DOI : 10.1007/s10928-017-9541-1], <https://hal.inria.fr/hal-01665722>
- [6] L. LÊ, M. BERGE, A. TFAYLI, A. BAILLET-GUFFROY, P. PROGNON, A. DOWEK, E. CAUDRON. *Quantification of gemcitabine intravenous drugs by direct measurement in chemotherapy plastic bags using a hand-held Raman spectrometer*, in "Talanta", May 2019, vol. 196, pp. 376-380, <https://hal.archives-ouvertes.fr/hal-01970020>
- [7] L. M. M. LÊ, M. BERGE, A. TFAYLI, J. ZHOU, P. PROGNON, A. BAILLET-GUFFROY, E. CAUDRON. *Rapid discrimination and quantification analysis of five antineoplastic drugs in aqueous solutions using Raman spectroscopy*, in "European Journal of Pharmaceutical Sciences", January 2018, vol. 111, pp. 158-166, <https://hal.archives-ouvertes.fr/hal-01970012>
- [8] L. M. M. LÊ, B. KÉGL, A. GRAMFORT, C. MARINI, D. NGUYEN, M. CHERTI, S. TFAYLI, A. TFAYLI, A. BAILLET-GUFFROY, P. PROGNON, P. CHAMINADE, E. CAUDRON. *Optimization of classification and regression analysis of four monoclonal antibodies from Raman spectra using collaborative machine learning approach*, in "Talanta", July 2018, vol. 184, pp. 260-265, <https://hal.archives-ouvertes.fr/hal-01969999>

- [9] L. LÊ, A. TFAYLI, P. PROGNON, E. CAUDRON. *Discriminative and Quantitative Analysis of Antineoplastic Taxane Drugs Using a Handheld Raman Spectrometer*, in "BioMed Research International ", July 2018, vol. 2018, pp. 1-7, <https://hal.archives-ouvertes.fr/hal-01969980>

International Conferences with Proceedings

- [10] N. BROSSE, A. DURMUS, É. MOULINES. *The promises and pitfalls of Stochastic Gradient Langevin Dynamics*, in "NeurIPS 2018 (Advances in Neural Information Processing Systems 2018)", Montreal, France, December 2018, <https://hal.inria.fr/hal-01934291>
- [11] G. ROBIN, H.-T. WAI, J. JOSSE, O. KLOPP, É. MOULINES. *Low-rank Interaction with Sparse Additive Effects Model for Large Data Frames*, in "32nd Conference on Neural Information Processing Systems (NeurIPS 2018)", Montréal, Canada, December 2018, <https://arxiv.org/abs/1812.08398> , <https://hal.archives-ouvertes.fr/hal-01959188>

Conferences without Proceedings

- [12] C. BIERNACKI, G. CELEUX, J. JOSSE, F. LAPORTE. *Dealing with missing data in model-based clustering through a MNAR model*, in "CMStatistics 2018 - 11th International Conference of the ERCIM WG on Computational and Methodological Statistics", Pise, Italy, December 2018, <https://hal.archives-ouvertes.fr/hal-01949120>

Other Publications

- [13] N. BROSSE, A. DURMUS, S. MEYN, É. MOULINES. *Diffusion approximations and control variates for MCMC*, November 2018, working paper or preprint, <https://hal.inria.fr/hal-01934316>
- [14] A. HAVET, M. LERASLE, É. MOULINES. *Density estimation for RWRE*, June 2018, <https://arxiv.org/abs/1806.05839> - working paper or preprint, <https://hal.archives-ouvertes.fr/hal-01815990>
- [15] F. HUSSON, J. JOSSE, B. NARASIMHAN, G. ROBIN. *Imputation of mixed data with multilevel singular value decomposition*, April 2018, <https://arxiv.org/abs/1804.11087> - working paper or preprint, <https://hal.archives-ouvertes.fr/hal-01781291>
- [16] W. JIANG, J. JOSSE, M. LAVIELLE. *Logistic Regression with Missing Covariates – Parameter Estimation, Model Selection and Prediction*, December 2018, working paper or preprint, <https://hal.archives-ouvertes.fr/hal-01958835>
- [17] B. KARIMI, M. LAVIELLE. *Efficient Metropolis-Hastings sampling for nonlinear mixed effects models*, December 2018, working paper or preprint, <https://hal.inria.fr/hal-01958247>
- [18] B. KARIMI, M. LAVIELLE, É. MOULINES. *f-SAEM: A fast Stochastic Approximation of the EM algorithm for nonlinear mixed effects models*, December 2018, working paper or preprint, <https://hal.inria.fr/hal-01958248>
- [19] L. MONTUELLE, E. LE PENNEC. *PAC-Bayesian aggregation of affine estimators*, January 2018, <https://arxiv.org/abs/1410.0661> - working paper or preprint, <https://hal.inria.fr/hal-01070805>
- [20] G. ROBIN, C. AMBROISE, S. S. ROBIN. *Incomplete graphical model inference via latent tree aggregation*, January 2018, working paper or preprint, <https://hal.archives-ouvertes.fr/hal-01686841>

- [21] A. SPORTISSE, C. BOYER, J. JOSSE. *Imputation and low-rank estimation with Missing Non At Random data*, January 2019, <https://arxiv.org/abs/1812.11409> - working paper or preprint, <https://hal.archives-ouvertes.fr/hal-01964720>