



**HAL**  
open science

## Learning spatio-temporal trajectories from manifold-valued longitudinal data

Jean-Baptiste Schiratti, Stéphanie Allasonniere, Olivier Colliot, Stanley Durrleman

► **To cite this version:**

Jean-Baptiste Schiratti, Stéphanie Allasonniere, Olivier Colliot, Stanley Durrleman. Learning spatio-temporal trajectories from manifold-valued longitudinal data. Neural Information Processing Systems, Dec 2015, Montréal, Canada. , Advances in Neural Information Processing Systems. hal-01245909

**HAL Id: hal-01245909**

**<https://hal.science/hal-01245909>**

Submitted on 18 Dec 2015

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

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

Public Domain

# Learning spatio-temporal trajectories from manifold-valued longitudinal data

Jean-Baptiste Schiratti<sup>1,2</sup>, Stéphanie Allasonnière<sup>2</sup>, Olivier Colliot<sup>1</sup>, Stanley Durrleman<sup>1</sup>

Inria Paris-Rocquencourt, Sorbonne Universités, UPMC Univ Paris 06 UMR S1127, CNRS UMR 7225, ICM, F-75013, Paris, France

CMAP, Ecole Polytechnique, Palaiseau, France

ARAMISLAB

www.aramislab.fr

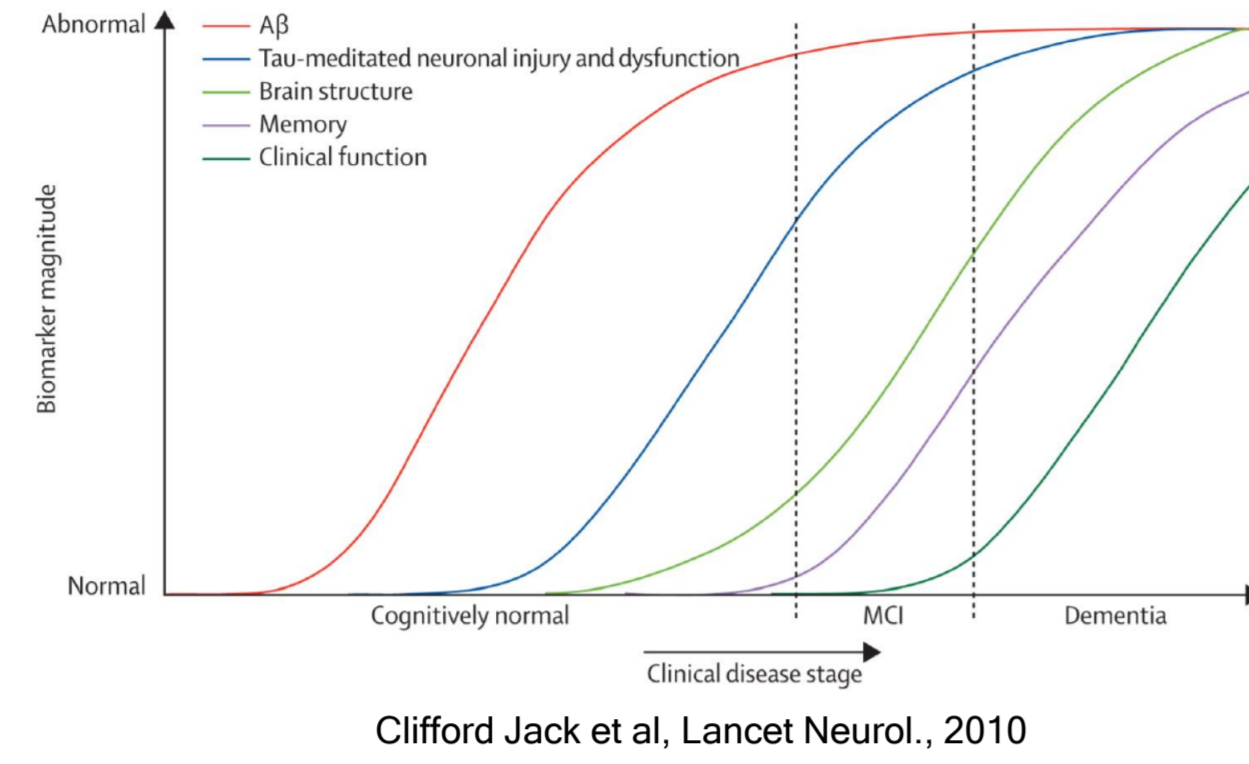


## Introduction

**Aim** : model the progression of neuro-degenerative diseases

Understanding the progression of neuro-degenerative diseases, such as Alzheimer's Disease (AD) is necessary for early and accurate diagnosis and care planning.

We need to validate experimentally hypothetical models of disease progression, such as [Clifford Jack et al, 2010].



Working with longitudinal data in the context of neuro-degenerative diseases raises two difficulties :

Two individuals of the same age might be at very different stages of disease progression

⇒ statistical models based on the regression of measurements with age are inadequate to model disease progression and age should not be treated as a covariate but as a random variable.

Longitudinal measurements sometimes belong to Riemannian manifolds (non-Euclidean spaces).

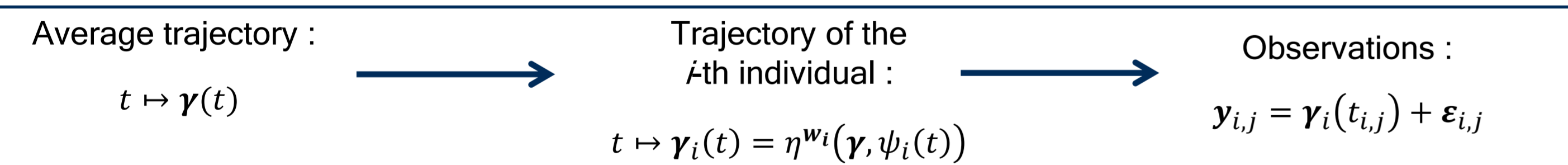
⇒ statistical models for such longitudinal data should be defined for manifold-valued measurements. Linear mixed-effects models [Laird and Ware, 1982] are not defined for manifold-valued measurements

## Generic spatio-temporal model for longitudinal data

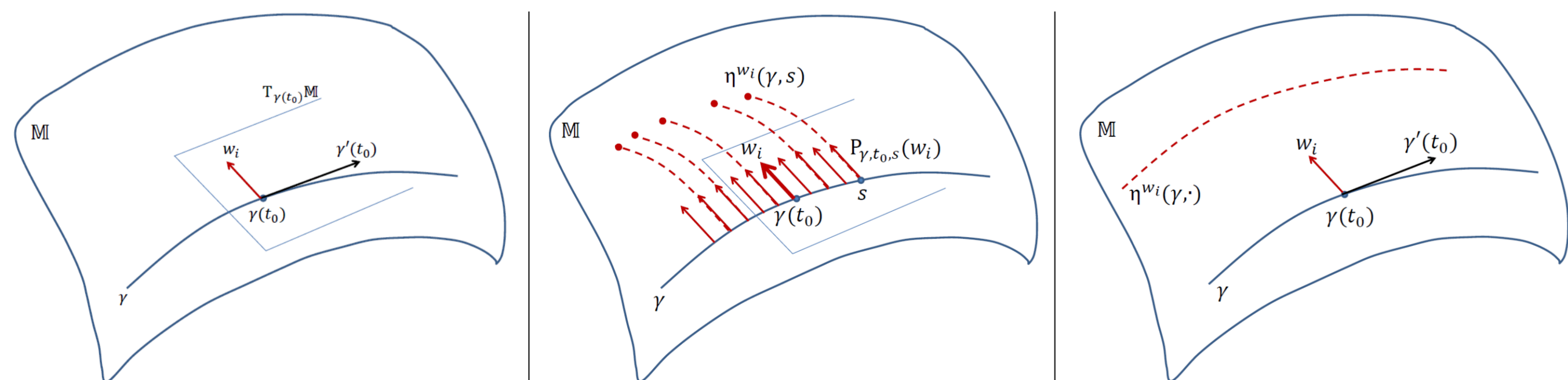
**Summary** : we propose a generic mixed-effects model for longitudinal manifold-valued data. The model allows to estimate an average trajectory as well as individual trajectories. Random effects allow to characterize changes in direction and pace at which individual trajectories are followed. This generic model is used to analyze the temporal progression of a family of univariate biomarkers.

- $(\bar{M}, g^{\bar{M}})$  smooth Riemannian manifold included in  $\mathbf{R}^n$
- $(M, g^M)$  sub-Riemannian manifold of  $\bar{M}$ , assumed to be geodesically complete
- $p \in M, v \in T_p M, \text{Exp}_p^M(v)$  : Riemannian exponential in  $M$  at  $p$  of the tangent vector  $v$
- $\gamma : \mathbf{R} \rightarrow M$  : geodesic of  $M$
- $t, t_0 \in \mathbf{R}, P_{\gamma, t_0, t}(\cdot)$  : parallel transport in  $M$  along  $\gamma$  from  $\gamma(t_0)$  to  $\gamma(t)$ .
- $t \mapsto \text{Exp}_{p, t_0}^M(v)(t)$  ; geodesic of  $M$  which goes through  $p$  at time  $t_0$  with velocity  $v$ .

➤ A hierarchical model :



- The average trajectory  $t \mapsto \gamma(t)$  is chosen to be the geodesic  $t \mapsto \text{Exp}_{p_0, t_0}^M(v_0)(t)$ ,  $p_0 \in M, v_0 \in T_{p_0} M$
- The trajectory of the  $i$ -th individual is obtained in two steps. We start by constructing the parallel shift of the average trajectory by using a tangent vector  $w_i$ , which we choose orthogonal to  $v_0$ .



- The trajectory  $t \mapsto \gamma_i(t)$  is then obtained by reparametrizing in time the parallel shift  $\eta^{w_i}(\gamma, \cdot)$  using the affine time reparametrization  $\psi_i : t \mapsto \alpha_i(t - t_0 - \tau_i) + t_0$ . This allows to account for the variability in stages of disease progression across the population.

The operation of parallel shifting, on the manifold  $M$ , using a tangent vector, is defined as follows :

**Definition** :  $w \in T_{\gamma(t_0)} M, w \neq 0$ . The curve  $s \mapsto \eta^w(\gamma, s)$  defined by :

$$\eta^w(\gamma, s) = \text{Exp}_{\gamma(s)}(P_{\gamma, t_0, s}(w)), s \in \mathbf{R}.$$

is said to be the « parallel shift of  $\gamma$  » using  $w$ .

By virtue of the tubular neighborhood theorem [Hirsch M.W., 2012], parallel shifting defines a local spatio-temporal coordinate system.

➤ The model :

$$y_{i,j} = \eta^{w_i}(\gamma, \psi_i(t_{i,j})) + \epsilon_{i,j}.$$

where  $\psi_i(t) = \alpha_i(t - t_0 - \tau_i) + t_0$  and :

- $\alpha_i$  is a subject-specific acceleration factor,  $\alpha_i = \exp(\xi_i)$  and  $\xi_i \sim N(0, \sigma_\xi^2)$
- $\tau_i$  is a subject-specific time shift,  $\tau_i \sim N(0, \sigma_\tau^2)$
- $\epsilon_{i,j} \sim N(0, \sigma^2)$
- $w_i$  is a subject-specific space shift :

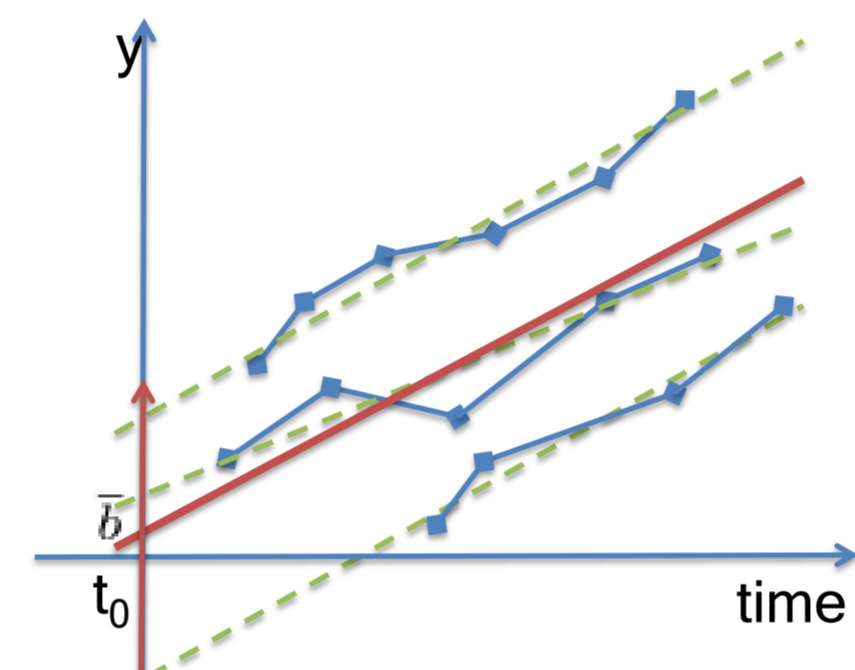
$$w_i = A s_i \text{ and } \forall 1 \leq j \leq N_s, s_{i,j} \sim \text{Laplace}\left(\frac{1}{2}\right).$$

In the spirit of Independent Component Analysis, the space shift  $w_i$  appears as a linear combination of  $N_s$  independent components, namely the columns of the matrix  $A$ .

## Three particular cases of our generic spatio-temporal model

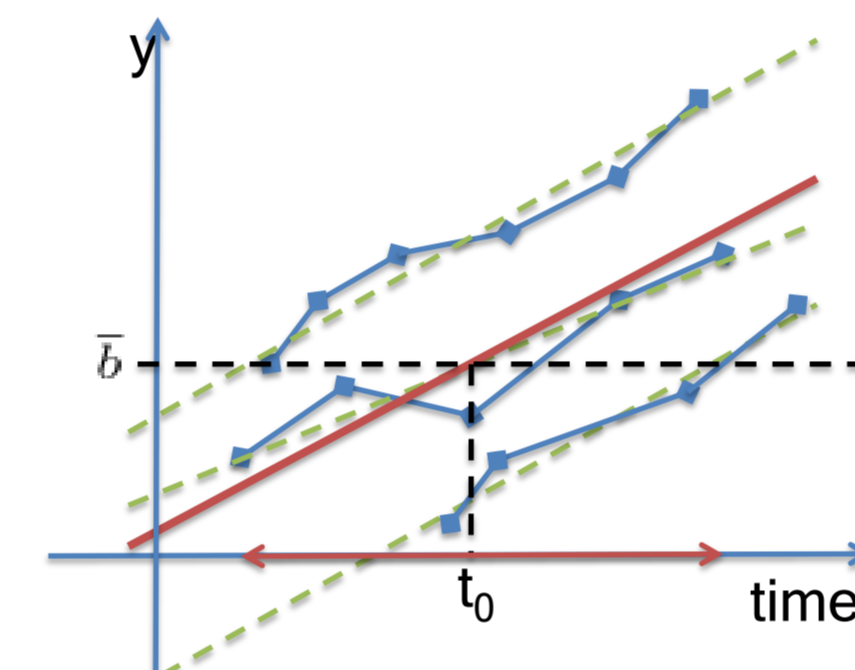
- « Straight lines » model [Schiratti et al., IPMI 2015]
- $M = \mathbf{R}$  (equipped with the canonical metric)
- Geodesics are straight lines

$$y_{i,j} = p_0 + \alpha_i v_0 (t_{i,j} - t_0 - \tau_i) + \epsilon_{i,j}.$$



$$y_{i,j} = (\bar{a} + \alpha_i)(t_{i,j} - t_0) + \bar{b} + b_i + \epsilon_{i,j}$$

[Laird & Ware, 1982]



$$y_{i,j} = (\bar{a} + \alpha_i)(t_{i,j} - t_0 - \tau_i) + \bar{b} + \epsilon_{i,j}$$

Our approach

- « Logistic curves » model [Schiratti et al., IPMI 2015]
- $M = ]0,1[$ , equipped with the metric  $g = (g_p)_{p \in ]0,1[}$ ,  $g_p(u, v) = uv / p^2(1-p)^2$
- Geodesics are logistic curves

$$y_{i,j} = \left(1 + \left(\frac{1}{p_0} - 1\right) \exp\left(-\frac{\alpha_i v_0 (t_{i,j} - t_0 - \tau_i)}{p_0(1-p_0)}\right)\right)^{-1} + \epsilon_{i,j}.$$

Note that this model is not equivalent to a linear model on the logit of the observations : the logit transform corresponds to the Riemannian logarithm at  $p_0 = 0.5$ . The model written in the tangent space is still *not* linear due to the multiplication between the random effects  $\alpha_i = \exp(\xi_i)$  and  $\tau_i$ .

• A progression model for a family of univariate biomarkers :

**Aim** : we want to analyze the temporal progression of a family of  $N$  biomarkers.

- We assume that the measurements of each biomarker belong to a one-dimensional Riemannian manifold  $I$ , geodesically complete and included in  $\mathbf{R}$ . As a consequence,  $M$  is a product of one-dimensional manifolds :  $M = I^N = I \times I \times \dots \times I$ .
- The average trajectory  $t \mapsto \gamma(t)$  is chosen among a parametric family of geodesics of  $M$  :

$$\gamma_\delta(t) = (\gamma_0(t), \gamma_0(t + \delta_1), \dots, \gamma_0(t + \delta_{N-1}))$$

where the parameters  $\delta_i$  ( $1 \leq i \leq N-1$ ) correspond to the relative delay between the biomarkers and  $t \mapsto \gamma_0(t)$  is a geodesic of the one-dimensional Riemannian manifold  $I$  (straight line, logistic curve, ...)

- If the space of observations is the open interval  $]0,1[$  for all the biomarkers (we consider normalized measurements), the manifold  $M$  is  $]0,1[^N$ , equipped with the product metric (see « Logistic curves model »).

Writing the generic spatio-temporal model in this case leads to a progression model for normalized biomarkers named « multivariate logistic curves model ». This model is given by :

$$y_{i,j,k} = \left(1 + \left(\frac{1}{p_0} - 1\right) \exp\left(-\frac{\alpha_i v_0 (t_{i,j} - t_0 - \tau_i) + v_0 \delta_k + v_0 \gamma'(t_0 + \delta_k)}{p_0(1-p_0)}\right)\right)^{-1} + \epsilon_{i,j,k}.$$

where  $y_{i,j,k}$  = measurement of the  $k$ -th biomarker for individual  $i$ , at time  $t_{i,j}$ .

## Estimation of the parameters of the model

The parameters of the generic spatio-temporal model are  $\theta = (p_0, t_0, v_0, \delta, \sigma_\xi, \sigma_\tau, \sigma, \text{vec}(A))$ .

**Summary** : the parameters are estimated using a stochastic version of the EM algorithm [Dempster, Laird, Rubin, 1977]. This algorithm is the Monte Carlo Markov Chain Stochastic Approximation EM algorithm (MCMC-SAEM) [Delyon et al., 1999 ; Allasonnière et al., 2010]. Theoretical results regarding the convergence of the algorithm have been proved in [Delyon et al., 1999 ; Allasonnière et al., 2010].

Note that the MCMC-SAEM requires that the model belongs to the *curved exponential family*. However, the multivariate logistic curves model does not belong to this family. The model can be made exponential by considering each parameters as realizations of independent Gaussian random variables.

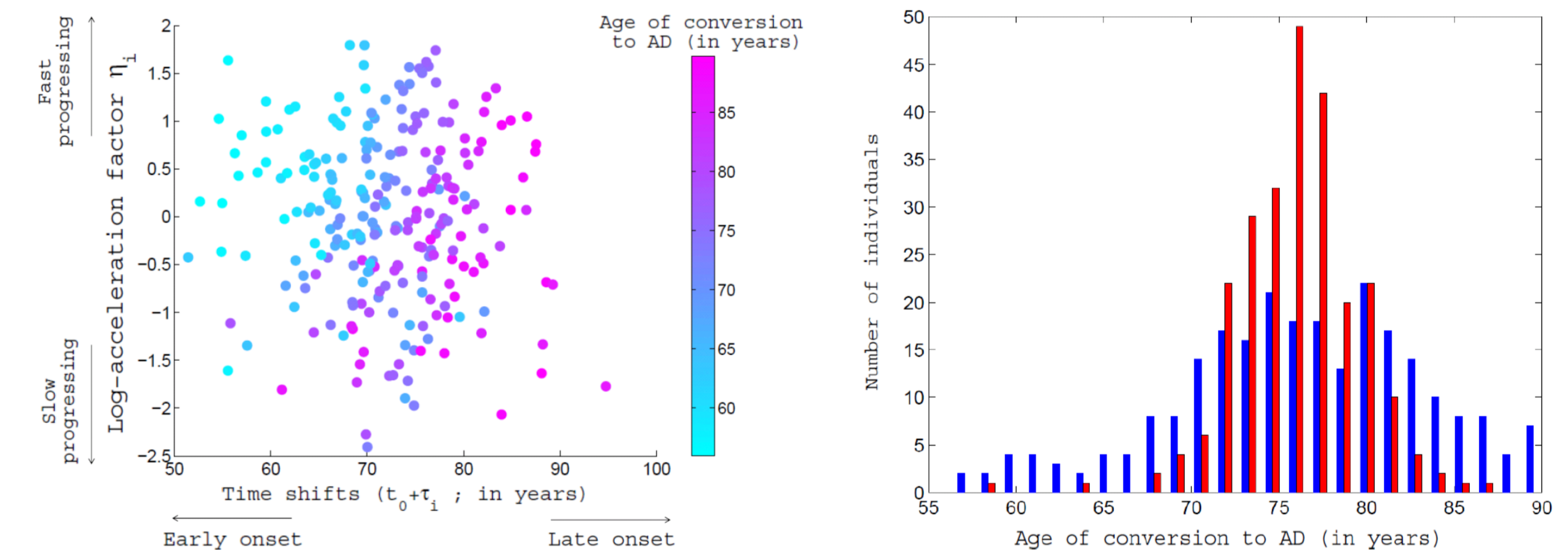
➤ Overview of the MCMC-SAEM for the multivariate logistic curves model :

$z^{(k)}$  (resp.  $\theta^{(k)}$ ) denotes the vector of hidden variables (resp. parameters) at the  $k$ -th iteration.

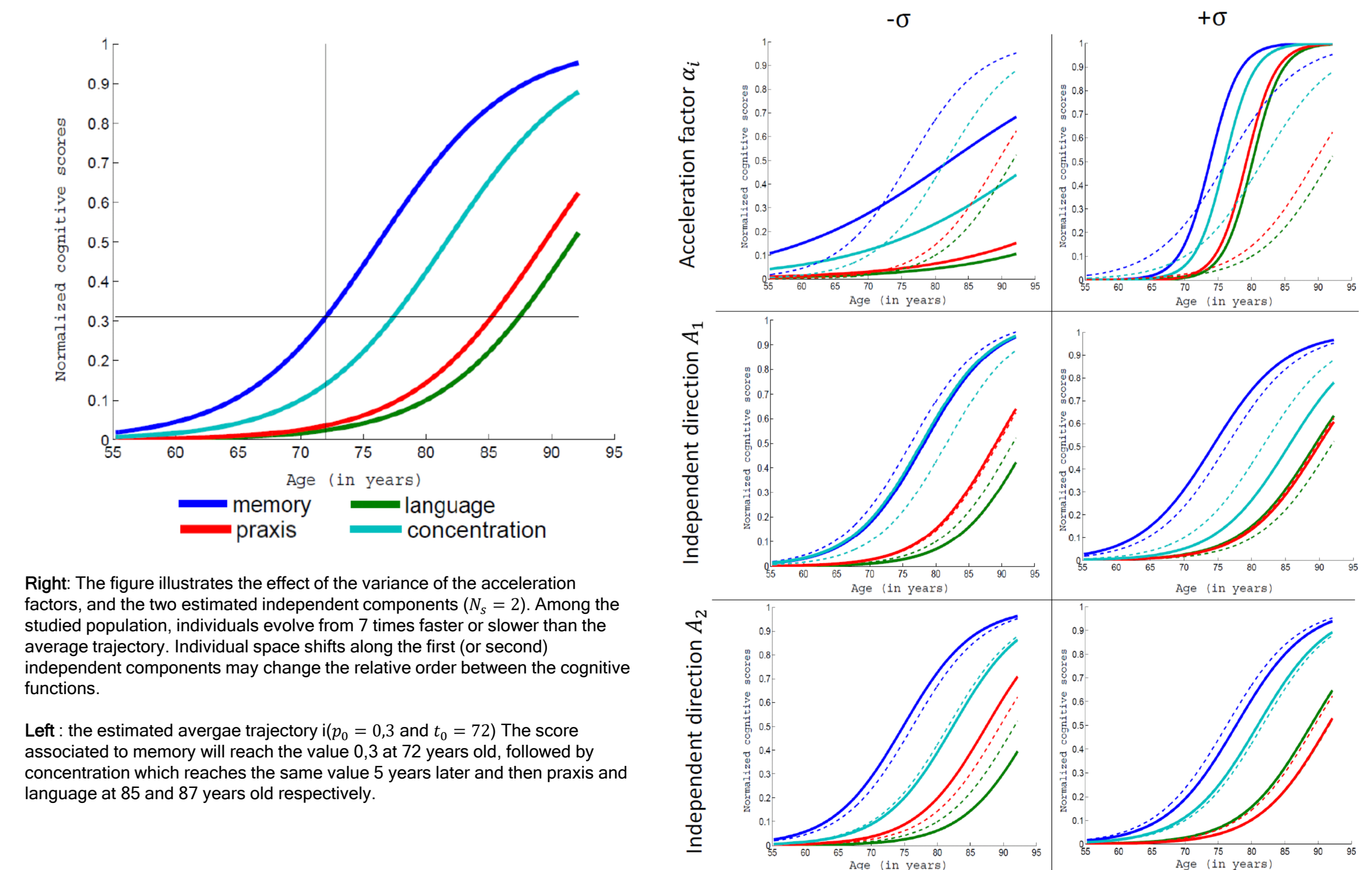
- Initialisation :
  - $\theta \leftarrow \theta^{(0)}, z^{(0)} \leftarrow \text{random}, S \leftarrow 0, (\epsilon_k)_k$  sequence of positive step-sizes
  - repeat until convergence
  - Simulation (Hasting-Metropolis within Gibbs sampler) :
    - $z^{(k+1)} \leftarrow \text{Gibbs sampler}(z^{(k)}, \theta^{(k)})$
    - Compute the sufficient statistics :  $S_1^{(k)} \leftarrow [y_{i,j}^T f_{i,j}]_{i,j}, S_2^{(k)} \leftarrow [\|f_{i,j}\|^2]_{i,j}, S_3^{(k)} \leftarrow [(\xi_i^{(k)})^2]_i,$
    - $S_4^{(k)} \leftarrow [(\tau_i^{(k)})^2]_i, S_5^{(k)} \leftarrow p_0^{(k)}, S_6^{(k)} \leftarrow t_0^{(k)}, S_7^{(k)} \leftarrow v_0^{(k)}, S_8^{(k)} \leftarrow [\delta_j^{(k)}]_j, S_9^{(k)} \leftarrow [\beta_j^{(k)}]_j$
    - Stochastic approximation :  $S_j^{(k+1)} \leftarrow S_j^{(k)} + \epsilon_k (S_j(z^{(k)}) - S_j^{(k)})$  for all  $j$
    - Maximization ( $\theta^{(k+1)} \leftarrow \text{argmax}_{\theta \in \Theta} [-\phi(\theta) + (S^{(k+1)}, \psi(\theta))]$ ) : closed-form updates
    - end repeat

## Experimental results

➤ **Data** : Normalized cognitive scores grouped into four categories (biomarkers) : memory (5 items), language (5 items), praxis (2 items), concentration (1 item). Data collected from the ADNI database for 248 MCI patients who converted to AD. Each observation is a point in  $M = ]0,1[^4$ .



Left : the plot of individual random effects show that the time-shifts correspond well with the age at which individuals converted to AD. Right : histogram of the ages of conversion to AD ( $t_i^{\text{conv}}$ ) in blue and histogram of the normalized ages ( $(\psi_i(t_i^{\text{conv}}))_i$ ) ; ages of conversion mapped from the individual timeline to the reference timeline using the subject-specific affine reparametrization. The individual time reparametrizations correctly register the dynamics of individual trajectories.



Right : The figure illustrates the effect of the variance of the acceleration factors, and the two estimated independent components ( $N_s = 2$ ). Among the studied population, individuals evolve from 7 times faster or slower than the average trajectory. Individual space shifts along the first (or second) independent components may change the relative order between the cognitive functions.

Left : the estimated average trajectory ( $p_0 = 0.3$  and  $t_0 = 72$ ) The score associated to memory will reach the value 0.3 at 72 years old, followed by concentration which reaches the same value 5 years later and then praxis and language at 85 and 87 years old respectively.