



**HAL**  
open science

# Monotonic Gaussian Process for Spatio-Temporal Trajectory Separation in Brain Imaging Data

Clement Abi Nader, Nicholas Ayache, Philippe Robert, Marco Lorenzi

► **To cite this version:**

Clement Abi Nader, Nicholas Ayache, Philippe Robert, Marco Lorenzi. Monotonic Gaussian Process for Spatio-Temporal Trajectory Separation in Brain Imaging Data. 2019. hal-02051843v1

**HAL Id: hal-02051843**

**<https://hal.science/hal-02051843v1>**

Preprint submitted on 28 Feb 2019 (v1), last revised 10 Oct 2019 (v3)

**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.

# Monotonic Gaussian Process for Spatio-Temporal Trajectory Separation in Brain Imaging Data

Clement Abi Nader<sup>1</sup>, Nicholas Ayache<sup>1</sup>, Philippe Robert<sup>2</sup>, and Marco Lorenzi<sup>1</sup>,  
for the Alzheimer’s Disease Neuroimaging Initiative\*

<sup>1</sup> UCA, Inria Sophia Antipolis, Epione Research Project

<sup>2</sup> UCA, CoBTeK lab, MNC3 program

**Abstract.** We introduce a probabilistic generative model for disentangling spatio-temporal disease trajectories from series of high-dimensional brain images. The model is based on spatio-temporal matrix factorization, where inference on the sources is constrained by anatomically plausible statistical priors. To model realistic trajectories, the temporal sources are defined as monotonic and time-reparametrized Gaussian Processes. To account for the non-stationarity of brain images, we model the spatial sources as sparse codes convolved at multiple scales. The method was tested on synthetic data favourably comparing with standard blind source separation approaches. The application on large-scale imaging data from a clinical study allows to disentangle differential temporal progression patterns mapping brain regions key to neurodegeneration, while revealing a disease-specific time scale associated to the clinical diagnosis.

## 1 Introduction

Neurodegenerative disorders such as Alzheimer’s disease (AD) are characterized by morphological and molecular changes of the brain, ultimately leading to cognitive and behavioral decline. Clinicians suggested hypothetical models of the disease evolution, showing how different types of biomarkers interact and lead to the final dementia stage [11]. In the past years, efforts have been made in order to collect large databases of imaging and clinical measures, hoping to obtain more insights about the disease progression through data-driven models describing the trajectory of the disease over time. This kind of models would

---

\*Data used in preparation of this article were obtained from the Alzheimers Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: [http://adni.loni.usc.edu/wp-content/uploads/how\\_to\\_apply/ADNI\\_Acknowledgement\\_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf).

be of critical importance for understanding the pathological progression in large scale data, and would represent a valuable reference for improving the individual diagnosis. Within this context, we propose a spatio-temporal generative model of disease progression, aimed at disentangling and quantifying the independent dynamics of changes observed in time-series of volumetric structural brain images. Moreover, thanks to our model, we aim at automatically inferring the disease severity of a patient with respect to the estimated trajectory. Defining such a disease progression model raises a number of methodological challenges.

AD spreads over decades with a temporal mismatch between the onset of the disease and the moment where the clinical symptoms appear. Either age of diagnosis, or the chronological age, are therefore not suitable as a temporal reference to describe the disease progression in time. Moreover, as the follow-up of patients doesn't exceed a few years, the development of a model of long-term pathological changes requires to integrate cross-sectional data from different individuals, in order to consider a longer period of time. In virtue of the lack of a well defined temporal reference, observations from different individuals are characterized by large and unknown variability in the onset and speed of the disease. It is therefore necessary to account for a time reparameterization function, mapping each individuals' observations to a common temporal axis associated to the absolute disease trajectory [12,27]. This would allow to estimate an absolute time-reference related to the natural history of the pathology.

The analysis of structural imaging data, such as the one provided by Magnetic Resonance Imaging (MRI), requires to account for spatio-temporally correlated features (voxels, i.e. volumetric pixels) defined over arrays of more than a million entries. The development of inference schemes jointly accounting for these correlation properties thus raises scalability issues, especially when accounting for the non-stationarity of the image signal. Furthermore the brain regions involved in AD exhibit various dynamics in time, and evolve at different speed [29]. From a modeling perspective, accounting for differential trajectories over space and time raises the problem of source identification and separation. This issue has been widely addressed in neuroimaging via Independent Component Analysis [6], especially on functional MRI (fMRI) data [5]. Nevertheless, while fMRI time-series are usually defined over a hundreds of time points acquired per subject, our problem consists in jointly analyzing short-term cross-sectional data observations with respect to an unknown time-line. This problem cannot be tackled with standard ICA, as time is generally an independent variable on which inference is not required. Moreover, ICA retrieves spatial sources based on an assumption of statistical independence. This assumption does not necessarily lead to clinically interpretable findings. Indeed, dependency across temporal patterns can be still highly relevant to the pathology, for example when modelling temporal delay across similar sources.

The problem of providing a realistic description of the biological processes is

critical when analyzing biomedical data, such as medical images. For example, to describe a plausible evolution of AD from normal to pathological stages the temporal sources need to be smooth and monotonic. It is also necessary to account for the non-stationarity of changes affecting the brain from global to localized spatio-temporal processes. As a result, spatial sources need to account for different resolutions at which these changes take place. While several multi-scale analysis approaches have been proposed to model spatio-temporal signals [19,4,10], extending this type of methods to the high-dimension of medical images is generally not trivial due to scalability issues. Finally, the noisy nature of medical images, along with the large signal variability across observations, requires a modelling framework robust to noise and over-fitting.

In this work, we propose to jointly address these issues within a Bayesian framework for spatio-temporal analysis of large-scale collections of volumetric medical images. We show that this framework allows us to naturally encode plausibility constraints through clinically-inspired priors, while accounting for the uncertainty of the temporal profiles and brain structures we wish to estimate. Similarly to the ICA setting, we formulate the problem of trajectory modelling as a matrix factorization across temporal and spatial sources. To promote smoothness in time and avoid any unnecessary hypothesis on the temporal trajectories, we rely on non-parametric modelling based on Gaussian Process (GP). We account for a plausible evolution from healthy to pathological stages thanks to a monotonicity constraint applied on the GP. Moreover, individuals' observations are temporally re-aligned on a common scale via a time-warping function. To model the non-stationarity of the spatial signal, the spatial sources are defined as sparse activation maps convolved at different scales. We show that our framework can be efficiently optimized through stochastic variational inference, allowing to exploit automatic differentiation and GPU support to speed up computations.

The paper is organized as follows: Section 2 analyzes related work on spatio-temporal modelling of neurodegeneration, while Section 3 details our method. In Section 4 we present experiments on synthetic data in which we compare our model to standard blind source separation approaches. We finally provide a demonstration of our method on the modelling of imaging data from a large scale clinical study. Prospects for future work and conclusions are drawn in section 5. Derivations that we could not fit in the paper are detailed in the Supplementary Material.

## 2 Related Work in Neurodegeneration Modelling

To deal with the uncertainty of the time-line of neurodegenerative pathologies, the concept of time-reparameterization of imaging-derived features has been used in several works. The underlying principle consists in estimating an absolute time-scale of disease progression by temporally re-aligning data from different subjects. For instance, in [30] the time-evolution was approximated as a sequence of events

which need to be re-ordered for each patient. This approach thus considers the evolution of neurodegenerative diseases as a collection of transitions between different clinical stages. This hypothesis is however limiting, as it doesn't reflect the continuity of changes affecting the brain along the course of the pathology.

To address this limitation, we rely on a continuous parameterization of the time-axis as in [18,8]. In particular, individuals' observations are time-realigned on a common temporal scale via a time-warping function. Using a set of relevant scalar biomarkers, these approaches allow to learn a time-scale describing the pathology evolution, and to estimate a time-line markedly correlated with the decline of cognitive abilities. Similarly, in [3] a disease progression score was estimated using biomarkers from molecular imaging. These methods are however based on the analysis of low-dimensional measures, such as collections of clinical variables. Therefore, they do not allow to scale to the high dimension of medical images. Our work tackles this shortcoming thanks to a scalable inference scheme based on stochastic variational inference.

Concerning the spatio-temporal representation of neurodegeneration, a mixed-effect model was proposed by [15] to learn an average spatio-temporal trajectory of brain evolution on cortical thickness data. The fixed-effect describes the average trajectory, while random effects are estimated through individual spatio-temporal warping functions, modelling how each subject differs from the global progression. Still, the extension of this approach to image volumes raises scalability issues. It has also to be noted that, to allow computational tractability, the brain evolution was assumed to be stationary in both space and time, thus limiting the ability of the model to disentangle the multiple dynamics of the brain structures involved in AD.

An attempt to sources separation is proposed in [20], through the decomposition of cortical thickness measurements as a mixture of spatio-temporal processes. This is performed by associating to each cortical vertex a temporal progression modeled by a sigmoid function, which may be however too simplistic to describe the progression of AD temporal processes. We propose to overcome this issue by non-parametric modelling of the temporal sources through GPs. Moreover, due to the lack of an explicit spatial correlation model, the approach of [20] may be potentially sensitive to spatial variation and noise, thus leading to poor interpretability. We address this problem by modelling the spatial sources through convolution of sparse maps at multiple resolutions, allowing to deal with signal non-stationarity and robustness to noise.

## 3 Methods

### 3.1 Individual time-shift

To account for the uncertainty of the time-line of individual measurements, we assume that the observations are defined with respect to an absolute temporal

reference  $\tau$ . This is performed through a time-warping function  $t_i = f^i(\tau)$ , that models the individual time-reparameterization. We choose a linear parameterization such that:

$$f^i(\tau) = \tau + \delta_i. \quad (1)$$

Within this setting the individual time-shift  $\delta_i$  encodes the temporal position of sample  $i$ , which in our application can be interpreted as the disease stage of subject  $i$  with respect to the long-term disease trajectory.

### 3.2 Data modelling

We assume that the spatio-temporal data is represented in a matrix  $\mathbf{Y}(x, f(\tau)) = [\mathbf{Y}_1(x, f^1(\tau)), \mathbf{Y}_2(x, f^2(\tau)), \dots, \mathbf{Y}_P(x, f^P(\tau))]^t$  with dimensions  $P \times F$ , where  $P$  is the number of samples,  $F$  the number of image features, and  $\mathbf{Y}_i(x, f^i(\tau))$  is a sample observed at position  $x$  and at time  $f^i(\tau)$ . We postulate a generative model in order to decompose the data in  $N_s$  spatio-temporal sources such that:

$$\mathbf{Y}_p(x, f^p(\tau)) = \mathbf{S}_p(\theta, f^p(\tau))\mathbf{A}(\psi, x) + \mathbf{Z} + \mathcal{E}. \quad (2)$$

Where  $\mathbf{S}$  is a  $P \times N_s$  matrix in which each column represents a temporal trajectory, and  $\theta$  the set of parameters related to the temporal sources.  $\mathbf{A}$  is a  $N_s \times F$  matrix where each row represents a spatial map, and  $\psi$  is the associated set of spatial parameters.  $\mathbf{Z}$  is a vector of length  $F$  that we need to estimate, representing structures which don't exhibit any intensity changes over time. Finally,  $\mathcal{E}$  follows a Gaussian distribution  $\mathcal{N}(0, \sigma^2 I)$ . According to the generative model, the data likelihood is:

$$p(\mathbf{Y}|\mathbf{A}, \mathbf{S}, \mathbf{Z}, \sigma) = \prod_{p=1}^P \frac{1}{(2\pi\sigma^2)^{\frac{F}{2}}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{Y}_p - \mathbf{S}_p\mathbf{A} - \mathbf{Z}\|^2\right). \quad (3)$$

Within a Bayesian modelling framework, we wish to maximize the marginal log-likelihood  $\log(p(\mathbf{Y}|\mathbf{Z}, \sigma))$ , to obtain posterior distributions for the spatio-temporal processes. Since the derivation of this quantity in a closed-form is not possible, we tackle this optimization problem through stochastic variational inference. Based on this formulation, in what follows we illustrate our model by detailing the variational approximations imposed on the spatio-temporal sources, along with the priors and constraints we impose to represent the data (Sections 3.3 and 3.4). Finally, we detail the variational lower bound and optimization strategy in Section 3.5.

### 3.3 Spatio-temporal processes

**Temporal sources.** In order to flexibly account for non-linear temporal patterns, the temporal sources are encoded in a matrix  $\mathbf{S}$  whose each column  $\mathbf{S}_n$  is a GP. To allow computational tractability within a variational setting, we rely on the GP approximation proposed in [7], through kernel approximation via random

feature expansion [25]. Within this framework, a GP can be approximated as a Bayesian Neural Network with form:  $\mathbf{S}_n(t) = \phi(\mathbf{\Omega}_n t) \mathbf{W}_n$ . For example, in the case of the Radial Basis Function (RBF) covariance,  $\mathbf{\Omega}_n$  is a linear projection in the spectral domain. It is equipped with a Gaussian distributed prior  $p(\mathbf{\Omega}_n) \sim \mathcal{N}(\mathbf{0}, l_n \mathbf{I})$  with a zero-mean and a covariance parameterized by a scalar  $l_n$ , acting as the length-scale parameter of the RBF covariance. The non-linear basis functions activation is defined by setting  $\phi(\cdot) = (\cos(\cdot), \sin(\cdot))$ , while the regression parameter  $\mathbf{W}_n$  is given with a standard normal prior. The GP inference problem can be conveniently performed by estimating approximated variational distributions for  $\mathbf{\Omega}_n$  and  $\mathbf{W}_n$  (Section 3.5).

We wish also to account for a steady evolution of the temporal processes, hence constraining the temporal sources to monotonicity. This is relevant in the medical case, where one would like to model the steady progression of a disease from normal to pathological stages. To do so, we constrain the space of the temporal sources to the set of solutions  $\mathcal{C}_n = \{\mathbf{S}_n(t) \mid \mathbf{S}'_n(t) \leq 0 \quad \forall t\}$ . This can be done consistently with the GP random feature expansion as shown in [17], where the constraint is introduced as a second likelihood term on the temporal sources dynamics:

$$p(\mathcal{C} \mid \mathbf{S}', \gamma) = (1 + \exp(-\gamma \mathbf{S}'(t)))^{-1}, \quad (4)$$

where  $\gamma$  controls the magnitude of the monotonicity constraint, and  $\mathcal{C} = \bigcap_n \mathcal{C}_n$ . According to [17] this constraint can be specified through the parametric form for the derivative of each  $\mathbf{S}_n$ :

$$\mathbf{S}'_n(t) = \mathbf{\Omega}_n \phi'(\mathbf{\Omega}_n t) \mathbf{W}_n. \quad (5)$$

This setting leads to an efficient scheme for estimating the temporal sources through stochastic variational inference (Section 3.5).

**Spatial sources.** According to the model introduced in Section 3.2, each observation  $\mathbf{Y}_p$  is obtained as the linear combination at a specific time-point between the temporal and spatial sources. In order to deal with the multi-scale nature of the imaging signal, we propose to represent the spatial sources at multiple resolutions. To this end, we encode the spatial sources in a matrix  $\mathbf{A}$  whose rows  $\mathbf{A}_n$  represent a specific source at a given scale. The scale is prescribed by a convolution operator  $\mathbf{\Sigma}_n$ , which is applied to a map  $\mathbf{B}_n$  that we wish to infer. This problem can be specified by defining  $\mathbf{A}_n = \mathbf{\Sigma}_n \mathbf{B}_n$ , where  $\mathbf{\Sigma}_n$  is an  $F \times F$  Gaussian kernel matrix imposing a specific spatial resolution. The length-scale parameter  $\lambda_n$  of the Gaussian kernel is fixed for each source, to force the model to pick details at that specific scale. Due to the high-dimension of the data we are modelling, performing stochastic variational inference in this setting raises scalability issues. For instance, if we assume a Gaussian distribution  $\mathcal{N}(\mu_{\mathbf{B}_n}, \text{diag}(\mathbf{\Lambda}))$  for  $\mathbf{B}_n$ , the distribution of the spatial signal would be  $p(\mathbf{A}_n) \sim \mathcal{N}(\mathbf{\Sigma}_n \mu_{\mathbf{B}_n}, \mathbf{\Sigma}_n \text{diag}(\mathbf{\Lambda}) \mathbf{\Sigma}_n^t)$ . As a result, sampling from  $p(\mathbf{A}_n)$  is not computationally tractable due to the size of the covariance matrix, which

prevents the use of standard inference schemes on  $\mathbf{B}_n$ . This can be overcome thanks to the separability of the Gaussian convolution kernel [21,16], according to which the 3D convolution matrix  $\boldsymbol{\Sigma}$  can be decomposed into the Kronecker product of 1D matrices,  $\boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}_n^x \otimes \boldsymbol{\Sigma}_n^y \otimes \boldsymbol{\Sigma}_n^z$ . This decomposition allows to efficiently perform standard operations such as matrix inversion, or matrix-vector multiplication [26]. Thanks to this choice, we recover tractability for the inference of  $\mathbf{B}_n$  through sampling, as required by stochastic inference methods [14].

### 3.4 Sparsity

In order to detect specific brain areas involved in neurodegeneration, we propose to introduce a sparsity constraint over the maps  $\mathbf{B}_n$ . Consistently with our variational inference scheme, we induce sparsity via *Variational Dropout* as proposed in [13]. This approach leverages on an improper log-scale uniform prior  $p(|\mathbf{B}_n|) \propto 1/|\mathbf{B}_n|$ , along with an approximate posterior distribution:

$$q_1(\mathbf{B}) = \prod_{n=1}^{Ns} \mathcal{N}(\boldsymbol{\mu}_n, \text{diag}(\alpha_{n,1}\boldsymbol{\mu}_{n,1}^2 \dots \alpha_{n,F}\boldsymbol{\mu}_{n,F}^2)). \quad (6)$$

In this formulation,  $\alpha_{n,f}$  is related to the individual dropout probability  $p_{n,f}$  of each weight by  $\alpha_{n,f} = p_{n,f}(1 - p_{n,f})^{-1}$ . To tackle the known stability issues affecting the inference of the dropout parameters, in what follows we leverage on the extension of *Variational Dropout* proposed in [22]. In this setting, the variance parameter is encoded in a new independent variable  $\boldsymbol{\rho}_{n,f}^2 = \alpha_{n,f}\boldsymbol{\mu}_{n,f}^2$ , while the posterior distribution is optimized with respect to  $(\boldsymbol{\mu}, \boldsymbol{\rho}^2)$ . Therefore, in order to minimize the cost function for large variance  $\boldsymbol{\rho}_{n,f}^2 \rightarrow \infty$  ( $\alpha_{n,f} \rightarrow \infty$  i.e.  $p_{n,f} \rightarrow 1$ ), the value of the weight's magnitude must be controlled by setting to zero the corresponding parameter  $\boldsymbol{\mu}_{n,f}$ . As a result, by dropping out weights in the code, we sparsify the estimated spatial maps, thus better isolating relevant spatial sub-structures. We note that although the elements of the spatial maps  $\mathbf{B}_n$  are assumed to be Gaussian and i.i.d, spatial correlations in the images are obtained thanks to the convolution operation detailed in Section 3.3.

### 3.5 Variational inference

To infer the individual time-shift parameter, the sets of parameters  $\theta$  and  $\psi$ , as well as  $\mathbf{Z}$  and  $\sigma$ , we need to jointly optimize the data evidence according to priors and constraints:

$$\begin{aligned} \log(p(\mathbf{Y}, \mathcal{C} | \mathbf{Z}, \sigma, \gamma)) &= \log \left[ \int p(\mathbf{Y} | \mathbf{B}, \mathbf{S}, \mathbf{Z}, \sigma) p(\mathcal{C} | \mathbf{S}', \gamma) p(\mathbf{B}) p(\mathbf{S}, \mathbf{S}' | \gamma) d\mathbf{B} d\mathbf{S} d\mathbf{S}' \right] \\ &= \log \left[ \int p(\mathbf{Y} | \mathbf{B}, \mathbf{S}, \mathbf{Z}, \sigma) p(\mathcal{C} | \mathbf{S}', \gamma) p(\mathbf{B}) p(\mathbf{S}' | \mathbf{S}, \gamma) p(\mathbf{S}) d\mathbf{B} d\mathbf{S} d\mathbf{S}' \right]. \end{aligned} \quad (7)$$



By observing that  $\mathbf{S}'$  is completely identified by  $\mathbf{S}$ , formula (7) can be written as:

$$\log(p(\mathbf{Y}, \mathcal{C} | \mathbf{Z}, \sigma, \gamma)) = \log \left[ \int p(\mathbf{Y} | \mathbf{B}, \mathbf{S}, \mathbf{Z}, \sigma) p(\mathcal{C} | \mathbf{S}', \gamma) p(\mathbf{B}) p(\mathbf{S}) d\mathbf{B} d\mathbf{S} \right]. \quad (8)$$

Since this integral is intractable, we tackle the optimization of (8) via stochastic variational inference. Following [7] and [17] we introduce approximations,  $q_2(\boldsymbol{\Omega})$  and  $q_3(\mathbf{W})$  in addition to  $q_1(\mathbf{B})$ , to derive the lower bound (a detailed derivation is given in the Supplementary Material):

$$\begin{aligned} \log(p(\mathbf{Y}, \mathcal{C} | \mathbf{Z}, \sigma, \gamma)) &\geq \mathbb{E}_{\mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}} [\log(p(\mathbf{Y} | \mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}, \mathbf{Z}, \sigma))] + \mathbb{E}_{\boldsymbol{\Omega}, \mathbf{W}} [\log(p(\mathcal{C} | \boldsymbol{\Omega}, \mathbf{W}, \gamma))] \\ &\quad - \mathcal{D}[q_1(\mathbf{B}) || p(\mathbf{B})] - \mathcal{D}[q_2(\boldsymbol{\Omega}) || p(\boldsymbol{\Omega})] - \mathcal{D}[q_3(\mathbf{W}) || p(\mathbf{W})]. \end{aligned} \quad (9)$$

Where  $\mathcal{D}$  refers to the Kullback-Leibler (KL) divergence.

The choice of prior and approximate posterior distribution for the maps  $\mathbf{B}_n$  leads to a closed-form formula for the KL divergence detailed in [22]:

$$-\mathcal{D}[q_1(\mathbf{B}) || p(\mathbf{B})] = \sum_{n,f} k_1 h(k_2 + k_3 \log(\alpha_{n,f})) - 0.5 \log(1 + \alpha_{n,f}^{-1}) - k_1, \quad (10)$$

where  $h$  is the sigmoid function and  $k_1, k_2, k_3$  are given constants. The approximated distributions  $q_2(\boldsymbol{\Omega})$  and  $q_3(\mathbf{W})$  are factorized across GPs such that:

$$\begin{aligned} q_2(\boldsymbol{\Omega}) &= \prod_{n=1}^{N_s} q_2(\boldsymbol{\Omega}_n) = \prod_{n=1}^{N_s} \prod_{j=1}^{N_{rf}} \mathcal{N}(r_{n,j}, p_{n,j}^2), \\ q_3(\mathbf{W}) &= \prod_{n=1}^{N_s} q_3(\mathbf{W}_n) = \prod_{n=1}^{N_s} \prod_{j=1}^{N_{rf}} \mathcal{N}(m_{n,j}, s_{n,j}^2), \end{aligned} \quad (11)$$

where  $N_{rf}$  is the number of random features used for the projection in the spectral domain. Using the Gaussian priors and the approximations we introduced above, we obtain the closed-form formula for the KL divergence between two Gaussian distributions, leading to:

$$\begin{aligned} \mathcal{D}[q_2(\boldsymbol{\Omega}) || p(\boldsymbol{\Omega})] &= \frac{1}{2} \sum_{n,j} \mathbf{p}_{n,j}^2 l_n + \mathbf{r}_{n,j}^2 l_n - 1 - \log(\mathbf{p}_{n,j}^2 l_n), \\ \mathcal{D}[q_3(\mathbf{W}) || p(\mathbf{W})] &= \frac{1}{2} \sum_{n,j} \mathbf{s}_{n,j}^2 + \mathbf{m}_{n,j}^2 - 1 - \log(\mathbf{s}_{n,j}^2). \end{aligned} \quad (12)$$

Finally we optimize the individual time-shifts  $\{\delta_i\}_{i=0}^P$ ,  $\mathbf{Z}$ ,  $\sigma$  as well as the subsequent sets of parameters:

$$\begin{aligned} \theta &= \{\mathbf{m}_n, \mathbf{s}_n^2, \mathbf{r}_n, \mathbf{p}_n^2, l_n, n \in [1, N_s]\}, \\ \psi &= \{\boldsymbol{\mu}_n, \boldsymbol{\rho}_n^2, n \in [1, N_s]\}. \end{aligned} \quad (13)$$

Following [14] and using the reparameterization trick, we can efficiently sample from the approximated distributions  $q_1, q_2$  and  $q_3$  to compute the two expectation terms from (9). We chose to alternate the optimization between the spatio-temporal parameters and the time-shift. We empirically set  $\gamma$  to the minimum value that gives monotonic sources, and the dropout rate of each  $\mathbf{B}_{n,f}$  to 95% (i.e  $\alpha_{n,f} = 19$ ), while  $\sigma$  was optimized during training along with the time-shift. The model is implemented and trained using the Pytorch library [23]. In the following sections we will refer to our method as Monotonic Gaussian Process Analysis (MGPA).

## 4 Experiments and Results

In this section we first benchmark MGPA on synthetic data to demonstrate its reconstruction and separation properties while comparing it to standard sources separation methods. We finally apply our model on a large set of medical images from a publicly available clinical study, demonstrating the ability of our method to retrieve spatio-temporal processes relevant to AD, along with a time-scale describing the course of the disease.

### 4.1 Synthetic tests on spatio-temporal trajectory separation

We tested MGPA on synthetic data generated as a linear combination of temporal functions and 3D activation maps at prescribed resolutions. The goal was to assess the method’s ability to identify the spatio-temporal sources underlying the data. We benchmarked our method with respect to ICA, Non-Negative Matrix Factorization (NMF), and Principal Component Analysis (PCA), which were applied from the standard implementation provided in the Scikit-Learn library [24].

The benchmark was specified by defining a 10-folds validation setting, generating the data at each fold as a linear combination of temporal sources  $\tilde{\mathbf{S}}(t) = [\tilde{\mathbf{S}}_1(t), \tilde{\mathbf{S}}_2(t)]$ , and spatial maps  $\tilde{\mathbf{A}} = [\tilde{\mathbf{A}}_1, \tilde{\mathbf{A}}_2]^t$ . The data was defined as  $\mathbf{Y}_p = \tilde{\mathbf{S}}(t_p)\tilde{\mathbf{A}} + \mathcal{E}_p$  over 50 time points  $t_p$ , where  $t_p$  was uniformly distributed in the range  $[0, 0.7]$ , and  $\mathcal{E}_p \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . The temporal sources were specified as sigmoid functions  $\tilde{\mathbf{S}}_i(t) = 1/(1 + \exp(-t + \alpha_i))$ , while the spatial structures had dimensions  $(30 \times 30 \times 30)$  such that  $\tilde{\mathbf{A}}_i = \tilde{\mathbf{\Sigma}}_i \tilde{\mathbf{B}}_i$ . The  $\tilde{\mathbf{\Sigma}}_i$  were chosen as Gaussian convolution matrices with respective length-scale of  $\lambda = 2$  mm and  $\lambda = 1$  mm. The  $\tilde{\mathbf{B}}_i$  were randomly sampled sparse 3D maps.

**Variable selection.** We applied our method by specifying an over-complete set of six sources with respective spatial length-scale of  $\lambda = 2, 2, 1, 1, 0.5, 0.5$  mm. Figure 1 shows an example of the sparse maps obtained for a specific fold. The model prunes the signal for most of the maps, while retaining two sparse maps,  $\mathbf{B}_0$  and  $\mathbf{B}_4$ , whose length-scale are  $\lambda = 2$  mm and  $\lambda = 1$  mm, thus correctly estimating the right number of sources and their spatial resolution. As

it can be qualitatively observed in Figure 1, we notice that the estimated sparse code convolved with a Gaussian kernel matrix with  $\lambda = 1$  mm is closer to its ground truth than the one convolved with a length-scale  $\lambda = 2$  mm. According to our tests, sparse codes associated to high resolution details (low  $\lambda$ ) are indeed more identifiable. On the contrary, the identifiability of images obtained via a convolution operator with larger kernels (large  $\lambda$ ) is lower, since these maps can be equivalently obtained through the convolution of different sparse codes.

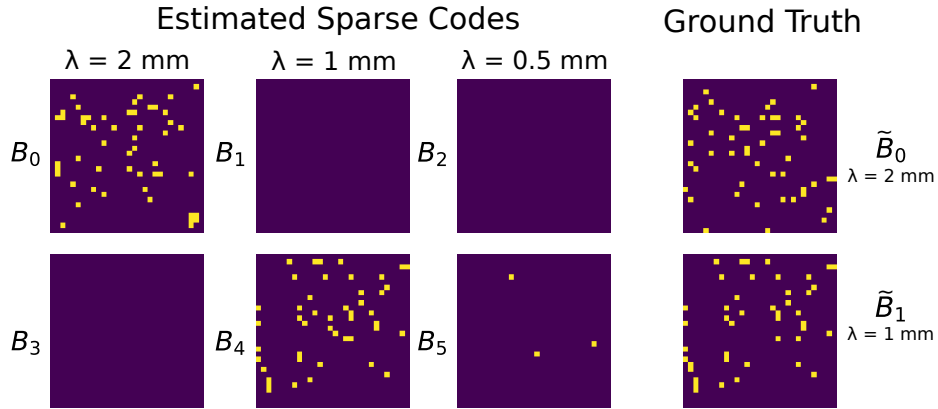


Fig. 1: Slices extracted from the six sparse codes and the ground truth. Blue: Rejected points. Yellow: Retained points.

**Sources separation.** We observe in Table 1 that the lowest Mean-Squared Error (MSE) for the temporal sources reconstruction is obtained by MGPA, closely followed by ICA. Similarly, our model and ICA show the highest Structural Similarity (SSIM) score [28], which quantifies the image reconstruction accuracy with respect to the ground truth maps, while accounting for the interdependencies between spatially close pixels. An example of image reconstruction from a sample fold is illustrated in Figure 2.

Table 1: MSE and SSIM between respectively the ground truth temporal and spatial sources with respect to the ones estimated by the different methods.

	TEMPORAL (MSE)	SPATIAL (SSIM)
MGPA	$(8 \pm 4) \cdot 10^{-5}$	$98\% \pm 1$
ICA	$(6 \pm 3) \cdot 10^{-4}$	$97\% \pm 2$
NMF	$(3 \pm 2) \cdot 10^{-2}$	$40\% \pm 17$
PCA	$0.44 \pm 10^{-3}$	$15\% \pm 1$

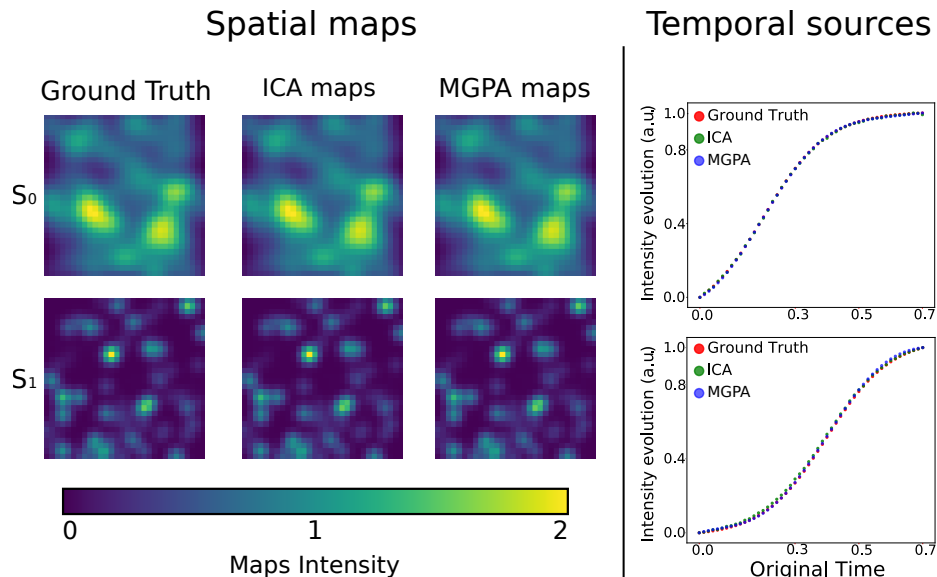


Fig. 2: Spatial maps: Sample slice from ground truth images ( $S_0$   $\lambda = 2$  mm,  $S_1$   $\lambda = 1$  mm), the maps estimated by ICA, and the ones estimated by MGPA. Temporal sources: Ground truth temporal sources (red) along with sources estimated by ICA (green) and MGPA (blue).

#### 4.2 Synthetic tests on trajectory separation and time-reparameterization

In this test, we modify the experimental benchmark by introducing a further element of variability associated to the time-axis. The temporal and spatial sources were modelled following the same procedure as in Section 4.1, however the observations were mixed along the temporal axis. To do so we generated longitudinal data as  $\mathbf{Y}_{p,j} = \tilde{\mathbf{S}}(t_p)\mathbf{A} + \mathbf{E}_j$ , by sampling between 1 and 10 images per time-point and randomly re-arranging them along the time-axis (cf. time-shift  $t_p$  of each observation at initialization in Figures 3 and 4, panel “Time-Shift”). As a result, the method needs to estimate the spatio-temporal sources, while reconstructing the original time-ordering. Since the model is agnostic of a time-scale, we note that the time-shift may have a different range than the original time-axis. However, its relative ordering should be consistent with the original time points. We fitted a linear regression model over the 10 folds between the original time and the estimated time-shift parameter, and obtained an average  $R^2$  coefficient of 0.98 with a standard deviation of 0.005 (cf. Tabel 2). This is illustrated for two different folds in the Time-Shift panel of Figures 3 and 4, where we observe a strong linear correlation with the original time-line, meaning that the algorithm correctly re-ordered the data with respect to the original

Table 2: MSE and SSIM between respectively the ground truth temporal and spatial sources with respect to the ones estimated by MGPA.  $R^2$  coefficient of the linear regression between the original time-line and the estimated time-shift.

	TEMPORAL (MSE)	SPATIAL (SSIM)	$R^2$
MGPA	$(2 \pm 0.8) \cdot 10^{-2}$	$95\% \pm 4$	$0.98 \pm 0.005$

time-axis. However, we notice in Table 2 that the MSE of the temporal sources significantly increased, due to the additional difficulty brought by the time-shift estimation. Indeed, in order to reconstruct the temporal signal we need to perfectly re-align hundreds of observations. This is the case in Figure 3 (optimal reconstruction result), where the time-shift correlates with the original time-line, allowing to distinguish every single observation and reconstruct the original temporal profiles. Whereas in Figure 4 (sub-optimal reconstruction result), the estimated time-shift doesn't exhibit a perfect fit, and generally underestimates the time-reparameterization for the later and earlier time points. This may be related to the challenging setting to reconstruct the time-line identified by the original temporal sources. Indeed, we observe that  $\mathbf{S}_0$  reaches a plateau for early time points, while  $\mathbf{S}_1$  is flat for later ones. This behaviour increases the difficulty of differentiating time points with low signal differences. As a result, it impacts the time-shift optimization and adds variability to the time-shift estimation performances, thus deteriorating the reconstruction of the temporal sources over the 10 folds compared to the previous benchmark. The spatial sources estimation remains comparable to the one without time-shift both quantitatively, with an average SSIM of 95%, and qualitatively, as shown in Figures 3 and 4. Within this setting, ICA, NMF and PCA poorly perform as they can't reconstruct the time-line. Interestingly, spatial ICA correctly estimated the spatial processes without however associating them to the corresponding temporal profile.

### 4.3 Application to spatio-temporal brain progression modelling

**Data.** Data used in the preparation of this article were obtained from the Alzheimers Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu).

We selected a cohort of 543 patients of the ADNI database composed of 88 controls (NL), 343 Mild Cognitive Impairment (MCI) and 118 AD patients at baseline. The MRI of each individual was processed following [1] in order to obtain gray matter (GM) density volumes in a standard anatomical space. These images have dimensions  $102 \times 130 \times 107$  before vectorization, leading to 1418820 spatial features per patient which represent the gray matter concentration of each voxel. From now on we will refer to the data as the  $(543 \times 1418820)$  matrix containing the images of all the subjects. We note that the analysis is performed by considering a single MRI scan per patient only. Therefore, the temporal

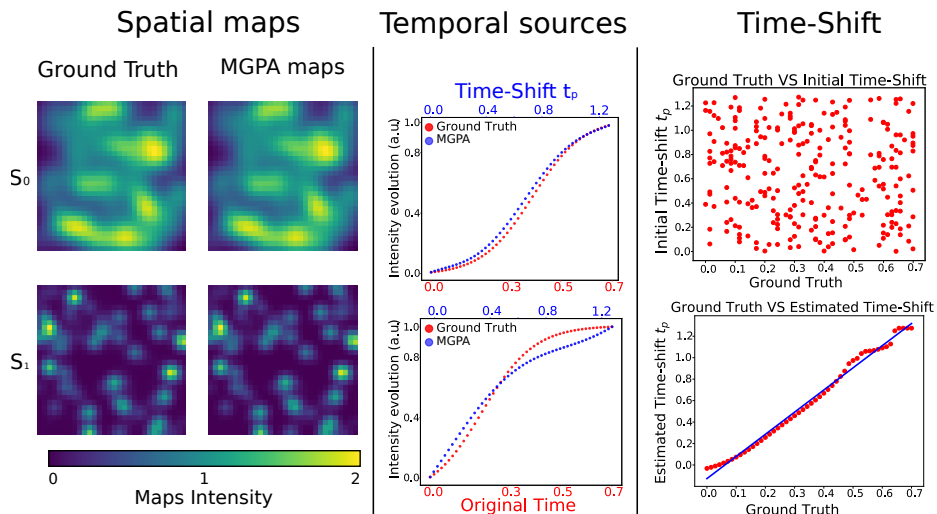


Fig. 3: Optimal reconstruction result on synthetic data. Spatial maps: Sample slice from ground truth images ( $S_0$   $\lambda = 2$  mm,  $S_1$   $\lambda = 1$  mm) and estimated spatial sources. Temporal sources: In red the original temporal sources, in blue the estimated temporal sources. Time-Shift: Time-shift  $t_p$  of each image at initialization (top), and after estimation (bottom). In blue, linear fit with the ground truth.

evolution has to be inferred solely through the analysis of relative differences between the brain morphologies across individuals.

**Model specification.** We aim at showing how MGPA applied on the MR images of the ADNI cohort is able to temporally re-align patients in order to describe AD progression in a plausible way, while detecting relevant spatio-temporal processes at stake in AD. To model the loss of gray matter over time the temporal sources are enforced to be monotonically decreasing. Since we don't consider any information about the disease stage of each individual before applying our method, all the observations are initialized at the same time reference  $\tau = 0$ . Therefore, as for the tests in Section 4.2, the time-shift reparameterization describes a relative re-ordering of the subjects not related to a specific time-unit. We apply our model by specifying an over-complete basis of six sources with  $\lambda = 1.5, 1.5, 0.75, 0.75, 0.1, 0.1$  mm, to cover both different scales and the variety of temporal evolution. Due to the high-dimension of the data matrix, the computations were parallelized over two GPUs, and the model required six hours to complete the training.

**Results.** In Figure 5 we show the four spatio-temporal processes retained by the model. The two sources with the highest resolution were discarded by our induced

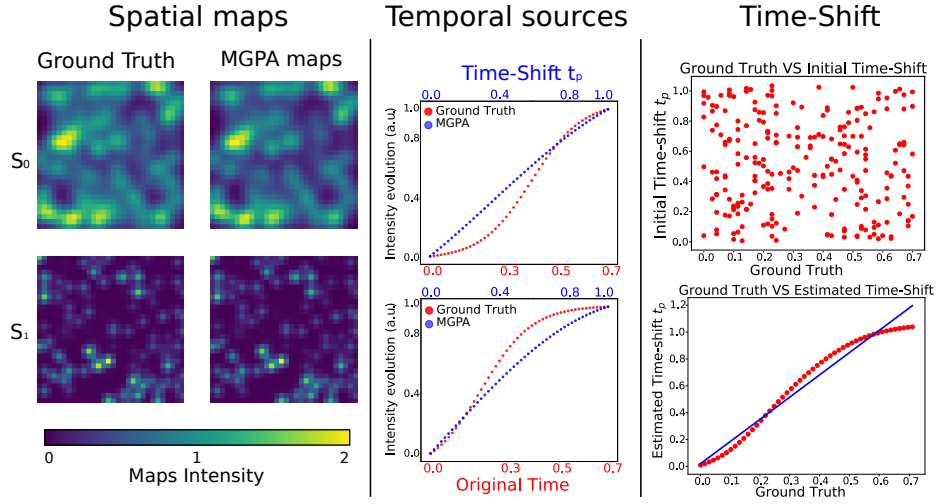


Fig. 4: Sub-optimal reconstruction result on synthetic data. Spatial maps: Sample slice from ground truth images ( $S_0$   $\lambda = 2$  mm,  $S_1$   $\lambda = 1$  mm) and estimated spatial sources. Temporal sources: In red the original temporal sources, in blue the estimated temporal sources. Time-Shift: Time-shift  $t_p$  of each image at initialization (top), and after estimation (bottom). In blue, linear fit with the ground truth.

sparsity, indicating the model robustness with respect to the noise from high-frequency signals. Sources  $S_0$  ( $\lambda = 1.5$  mm) and  $S_2$  ( $\lambda = 0.75$  mm) encompass a large extent of the brain with a focus on cortical areas, and exhibit a gray matter decrease that tends to plateau in the latest stages. It is also interesting to note the symmetry of the estimated spatio-temporal processes, showing similar accelerating progressions at two different spatial resolutions. Moreover, sources  $S_1$  ( $\lambda = 1.5$  mm) and  $S_3$  ( $\lambda = 0.75$  mm) exhibit a differential pattern of gray matter loss accelerating in the latest stages of the pathology. We note that the maps target subcortical areas such as the hippocampi, which are key regions of the AD pathology. These results underline the complex evolution of the brain gray matter, and the ability of the model to disentangle spatio-temporal processes mapping regions involved in the pathology [2,9].

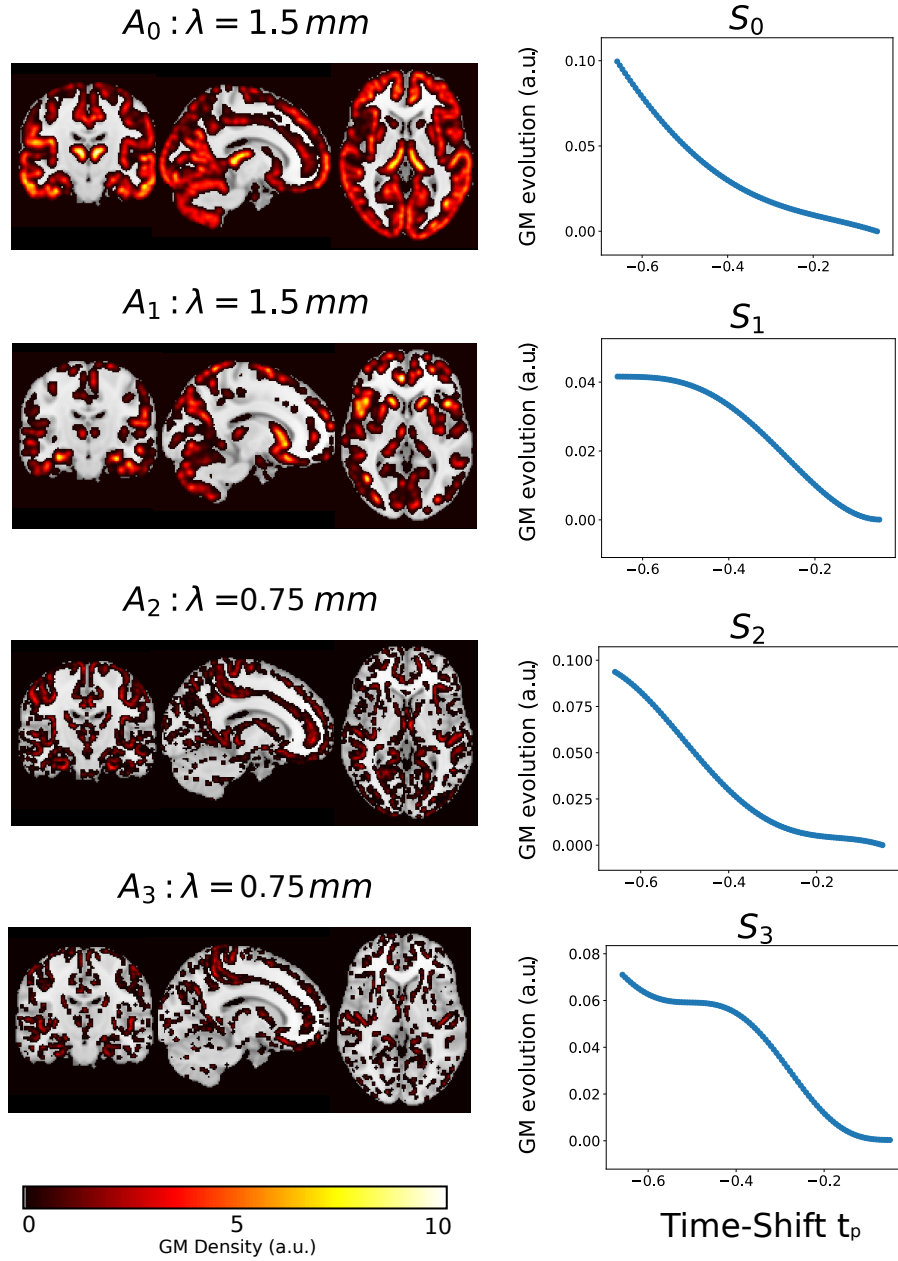


Fig. 5: The four estimated spatio-temporal processes underlying AD progression.



**Model Consistency.** To verify the plausibility of the estimated temporal reparameterization function, we compared the gray matter concentration of different brain regions against the time-shift value for each individual (Figure 6 top row). We observe a decrease of gray matter in brain regions as we progress along the estimated time-line, allowing to relate large time-shift values to lower gray matter density. This is confirmed by the agreement between the gray matter density predicted by the model and the raw concentration measures (Figure 6 bottom row).

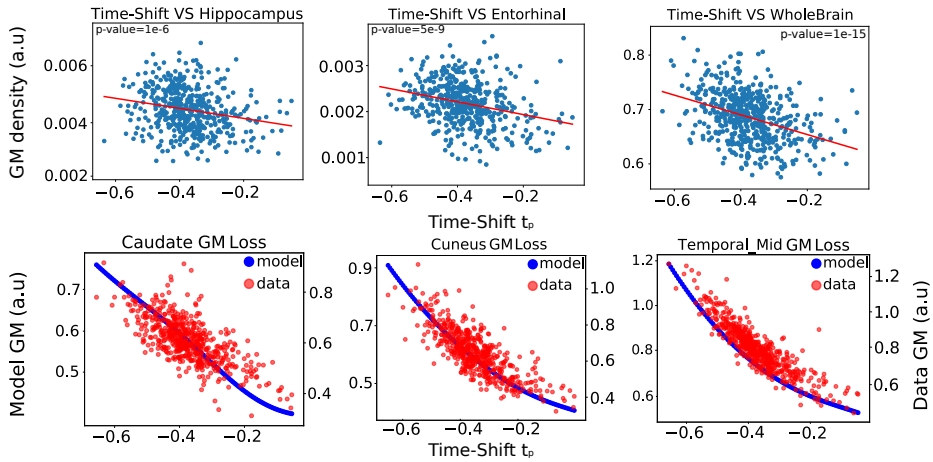


Fig. 6: Top row: Individuals’ volumetric biomarkers against time-shift  $t_p$ . Bottom row: Predicted gray matter density averaged on specific brain areas (blue line) and observed values (red dots), along the estimated time-line.

**Plausibility with respect to clinical evidence.** We assessed the clinical relevance of the estimated time-shift by relating it to independent medical information evaluated by physicians. To this end, we compared the estimated time-shift to the ADAS11 scale: high values of this score indicates a decline of cognitive abilities. We show in Figure 7 a non-linear relationship between ADAS11 and the time-shift, suggesting an acceleration of symptoms along the estimated time-course, which is characteristic of AD in its latest stages. The box-plot of Figure 7 shows the time-shift distribution across clinical groups. We observe an increase of the estimated time-shift when going from healthy to pathological stages. The high uncertainty associated to the MCI group is due to the broad definition of this clinical category, which includes subjects not necessarily affected by dementia. We note that the MCI subjects subsequently converted to AD (MCI to AD) exhibit higher time-shift than the MCI group, highlighting the ability of the model to differentiate clinical diagnosis without any prior knowledge. A

similar distinction can be noticed between the NL and NL to MCI groups. It is important to recall that this result is obtained from the analysis of a single MRI scan per patient only.

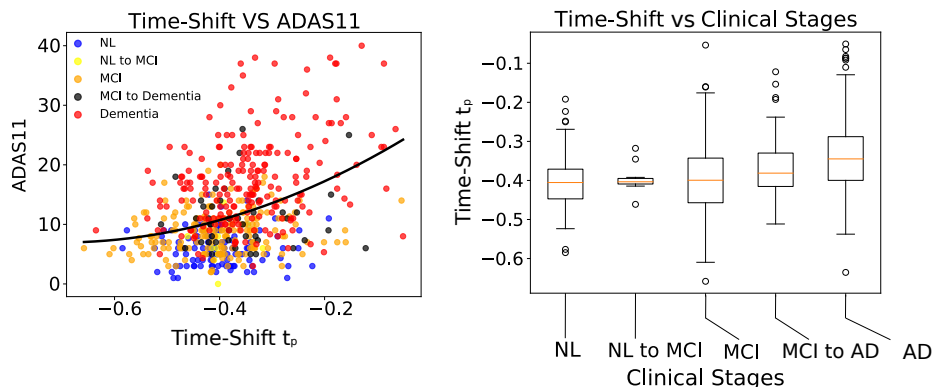


Fig. 7: Left: Evolution of the ADAS11 score along the estimated time-course. Right: Distribution of the time-shift values over the different clinical stages.

## 5 Discussion

We presented a generative model to analyze spatio-temporal data based on matrix factorization across temporal and spatial sources. The proposed application on large set of medical images show the ability of the model to disentangle relevant spatio-temporal processes at stake in AD, along with an estimated time-scale related to the disease evolution.

There are several avenues of improvement for the proposed approach. The optimization is highly sensitive to the initialization of the spatial sources. This is typical of such complex non-convex problems, and requires further investigations to better control the algorithm convergence. Moreover, as mentioned in Section 3.4, the *Variational Dropout* framework leads to stability issues affecting inference, which are mostly due to the improper prior. This problem motivates the need of alternative ways to induce sparsity on the spatial maps. The modelling results are also sensitive to the specification of the spatio-temporal processes priors. In our case, the monotonicity constraint imposed to the GPs may be too restrictive to completely capture the complexity of the progression of neurodegeneration. Ultimately, even if we essentially focused on the medical case, our approach remains general enough to be applied on different types of spatio-temporal data.

## References

1. Ashburner, J.: A fast diffeomorphic image registration algorithm. *NeuroImage* **38**(1), 95 – 113 (2007). <https://doi.org/https://doi.org/10.1016/j.neuroimage.2007.07.007>
2. Bateman, R.J., Xiong, C., Benzinger, T.L., Fagan, A.M., Goate, A., Fox, N.C., Marcus, D.S., Cairns, N.J., Xie, X., Blazey, T.M., Holtzman, D.M., Santacruz, A., Buckles, V., Oliver, A., Moulder, K., Aisen, P.S., Ghetti, B., Klunk, W.E., McDade, E., Martins, R.N., Masters, C.L., Mayeux, R., Ringman, J.M., Rossor, M.N., Schofield, P.R., Sperling, R.A., Salloway, S., Morris, J.C.: Clinical and biomarker changes in dominantly inherited alzheimer’s disease. *New England Journal of Medicine* **367**(9), 795–804 (2012), PMID: 22784036
3. Bilgel, M., Jedynak, B., Wong, D.F., Resnick, S.M., Prince, J.L.: Temporal Trajectory and Progression Score Estimation from Voxelwise Longitudinal Imaging Measures: Application to Amyloid Imaging. *Inf Process Med Imaging* **24**, 424–436 (2015)
4. Bullmore, E., Fadili, J., Maxim, V., Sendur, L., Whitcher, B., Suckling, J., Brammer, M., Breakspear, M.: Wavelets and functional magnetic resonance imaging of the human brain. *Neuroimage* **23 Suppl 1**, S234–249 (2004)
5. Calhoun, V.D., Liu, J., Adali, T.: A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data. *Neuroimage* **45**(1 Suppl), S163–172 (Mar 2009)
6. Comon, P.: Independent Component Analysis, a new concept? *Signal Processing* **36**, 287–314 (Apr 1994). [https://doi.org/10.1016/0165-1684\(94\)90029-9](https://doi.org/10.1016/0165-1684(94)90029-9)
7. Cutajar, K., Bonilla, E.V., Michiardi, P., Filippone, M.: Random feature expansions for deep Gaussian processes. In: Precup, D., Teh, Y.W. (eds.) *Proceedings of the 34th International Conference on Machine Learning. Proceedings of Machine Learning Research*, vol. 70, pp. 884–893. PMLR, International Convention Centre, Sydney, Australia (06–11 Aug 2017)
8. Donohue, M.C., Jacqmin-Gadda, H., Goff, M.L., Thomas, R.G., Raman, R., Gamst, A.C., Beckett, L.A., Jack, C.R., Weiner, M.W., Dartigues, J.F., Aisen, P.S.: Estimating long-term multivariate progression from short-term data. *Alzheimer’s & Dementia* **10**(5, Supplement), S400 – S410 (2014). <https://doi.org/https://doi.org/10.1016/j.jalz.2013.10.003>
9. Frisoni, G.B., Fox, N.C., Jack, C.R., Scheltens, P., Thompson, P.M.: The clinical use of structural MRI in Alzheimer disease. *Nat Rev Neurol* **6**(2), 67–77 (Feb 2010)
10. Hackmack, K., Paul, F., Weygandt, M., Allefeld, C., Haynes, J.D.: Multi-scale classification of disease using structural MRI and wavelet transform. *Neuroimage* **62**(1), 48–58 (Aug 2012)
11. Jack, C.R., Knopman, D.S., Jagust, W.J., Shaw, L.M., Aisen, P.S., Weiner, M.W., Petersen, R.C., Trojanowski, J.Q.: Hypothetical model of dynamic biomarkers of the Alzheimer’s pathological cascade. *Lancet Neurol* **9**(1), 119–128 (Jan 2010)
12. Jedynak, B.M., Lang, A., Liu, B., Katz, E., Zhang, Y., Wyman, B.T., Raunig, D., Jedynak, C.P., Caffo, B., Prince, J.L.: A computational neurodegenerative disease progression score: method and results with the Alzheimer’s disease Neuroimaging Initiative cohort. *Neuroimage* **63**(3), 1478–1486 (Nov 2012)
13. Kingma, D.P., Salimans, T., Welling, M.: Variational dropout and the local reparameterization trick. *CoRR* **abs/1506.02557** (2015)
14. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. *CoRR* **abs/1312.6114** (2013)

15. Koval, I., Schiratti, J.B., Routier, A., Bacci, M., Colliot, O., Allasonnière, S., Durrleman, S.: Statistical learning of spatiotemporal patterns from longitudinal manifold-valued networks. In: Medical Image Computing and Computer Assisted Intervention. Medical Image Computing and Computer Assisted Intervention, Quebec City, Canada (Sep 2017)
16. Lorenzi, M., Ziegler, G., Alexander, D.C., Ourselin, S.: Efficient Gaussian Process-Based Modelling and Prediction of Image Time Series. *Inf Process Med Imaging* **24**, 626–637 (2015)
17. Lorenzi, M., Filippone, M.: Constraining the dynamics of deep probabilistic models. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 3233–3242. PMLR, Stockholm, Sweden (10–15 Jul 2018)
18. Lorenzi, M., Filippone, M., Frisoni, G.B., Alexander, D.C., Ourselin, S.: Probabilistic disease progression modeling to characterize diagnostic uncertainty: Application to staging and prediction in alzheimer’s disease. *NeuroImage* (2017). <https://doi.org/https://doi.org/10.1016/j.neuroimage.2017.08.059>
19. Mallat, S.G.: A theory for multiresolution signal decomposition: The wavelet representation. *IEEE Trans. Pattern Anal. Mach. Intell.* **11**(7), 674–693 (Jul 1989). <https://doi.org/10.1109/34.192463>
20. Marinescu, R.V., Eshaghi, A., Lorenzi, M., Young, A.L., Oxtoby, N.P., Garbarino, S., Shakespeare, T.J., Crutch, S.J., Alexander, D.C.: A vertex clustering model for disease progression: Application to cortical thickness images. In: Niethammer, M., Styner, M., Aylward, S., Zhu, H., Oguz, I., Yap, P.T., Shen, D. (eds.) Information Processing in Medical Imaging. pp. 134–145. Springer International Publishing, Cham (2017)
21. Marquand, A.F., Brammer, M., Williams, S.C., Doyle, O.M.: Bayesian multi-task learning for decoding multi-subject neuroimaging data. *Neuroimage* **92**, 298–311 (May 2014)
22. Molchanov, D., Ashukha, A., Vetrov, D.: Variational dropout sparsifies deep neural networks. In: Precup, D., Teh, Y.W. (eds.) Proceedings of the 34th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 70, pp. 2498–2507. PMLR, International Convention Centre, Sydney, Australia (06–11 Aug 2017)
23. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
24. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
25. Rahimi, A., Recht, B.: Random features for large-scale kernel machines. In: Platt, J.C., Koller, D., Singer, Y., Roweis, S.T. (eds.) Advances in Neural Information Processing Systems 20, pp. 1177–1184. Curran Associates, Inc. (2008)
26. Saatçi, Y.: Scalable inference for structured gaussian process models (2011)
27. Schiratti, J., Allasonnière, S., Colliot, O., Durrleman, S.: Learning spatiotemporal trajectories from manifold-valued longitudinal data. In: NIPS. pp. 2404–2412 (2015)
28. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: From error visibility to structural similarity. *IEEE TRANSACTIONS ON IMAGE PROCESSING* **13**(4), 600–612 (2004)
29. Whitwell, J.L.: Progression of atrophy in Alzheimer’s disease and related disorders. *Neurotox Res* **18**(3-4), 339–346 (Nov 2010)

30. Young, A.L., Oxtoby, N.P., Huang, J., Marinescu, R.V., Daga, P., Cash, D.M., Fox, N.C., Ourselin, S., Schott, J.M., Alexander, D.C.: Multiple Orderings of Events in Disease Progression. *Inf Process Med Imaging* **24**, 711–722 (2015)

## Supplementary Material

### A. Derivation of the Lower Bound

We give a detailed derivation of the lower bound from equation (9) that we use in variational inference to learn the different parameters of our model.

$$\begin{aligned}
\log p(\mathbf{Y}, \mathcal{C} | \mathbf{Z}, \sigma, \gamma) &= \log \left[ \int p(\mathbf{Y} | \mathbf{B}, \mathbf{S}, \mathbf{Z}, \sigma) p(\mathcal{C} | \mathbf{S}', \gamma) p(\mathbf{B}) p(\mathbf{S}) d\mathbf{B} d\mathbf{S} \right] \\
&= \log \left[ \int p(\mathbf{Y} | \mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}, \mathbf{Z}, \sigma) p(\mathcal{C} | \boldsymbol{\Omega}, \mathbf{W}, \gamma) p(\mathbf{B}) p(\boldsymbol{\Omega}) \right. \\
&\quad \left. p(\mathbf{W}) d\mathbf{B} d\boldsymbol{\Omega} d\mathbf{W} \right] \\
&= \log \left[ \int p(\mathbf{Y} | \mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}, \mathbf{Z}, \sigma) p(\mathcal{C} | \boldsymbol{\Omega}, \mathbf{W}, \gamma) p(\mathbf{B}) p(\boldsymbol{\Omega}) \right. \\
&\quad \left. p(\mathbf{W}) \frac{q_1(\mathbf{B}) q_2(\boldsymbol{\Omega}) q_3(\mathbf{W})}{q_1(\mathbf{B}) q_2(\boldsymbol{\Omega}) q_3(\mathbf{W})} d\mathbf{B} d\boldsymbol{\Omega} d\mathbf{W} \right] \\
&= \log \left[ \mathbb{E}_{q_1, q_2, q_3} \frac{p(\mathbf{Y} | \mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}, \mathbf{Z}, \sigma) p(\mathcal{C} | \boldsymbol{\Omega}, \mathbf{W}, \gamma) p(\mathbf{B})}{q_1(\mathbf{B}) q_2(\boldsymbol{\Omega}) q_3(\mathbf{W})} \right. \\
&\quad \left. \frac{p(\boldsymbol{\Omega}) p(\mathbf{W})}{q_1(\mathbf{B}) q_2(\boldsymbol{\Omega}) q_3(\mathbf{W})} \right] \\
&\geq \mathbb{E}_{q_1, q_2, q_3} \left( \log \left[ \frac{p(\mathbf{Y} | \mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}, \mathbf{Z}, \sigma) p(\mathcal{C} | \boldsymbol{\Omega}, \mathbf{W}, \gamma) p(\mathbf{B})}{q_1(\mathbf{B}) q_2(\boldsymbol{\Omega}) q_3(\mathbf{W})} \right. \right. \\
&\quad \left. \left. \frac{p(\boldsymbol{\Omega}) p(\mathbf{W})}{q_1(\mathbf{B}) q_2(\boldsymbol{\Omega}) q_3(\mathbf{W})} \right] \right) \\
&= \mathbb{E}_{q_1, q_2, q_3} [\log(p(\mathbf{Y} | \mathbf{B}, \boldsymbol{\Omega}, \mathbf{W}, \mathbf{Z}, \sigma))] \\
&\quad + \mathbb{E}_{q_2, q_3} [\log(p(\mathcal{C} | \boldsymbol{\Omega}, \mathbf{W}, \gamma))] - \mathcal{D}[q_1(\mathbf{B}) || p(\mathbf{B})] \\
&\quad - \mathcal{D}[q_2(\boldsymbol{\Omega}) || p(\boldsymbol{\Omega})] - \mathcal{D}[q_3(\mathbf{W}) || p(\mathbf{W})].
\end{aligned}$$

### B. ADNI

Data collection and sharing for this project was funded by the Alzheimer’s Disease Neuroimaging Initiative (ADNI) and DOD ADNI. ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimers Association; Alzheimers Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech,

Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health ([www.fnih.org](http://www.fnih.org)). The grantee organization is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimers Therapeutic Research Institute at the University of Southern California. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of Southern California.