



A Multimodal Disease Progression Model for Genetic Associations with Disease Dynamics

Nemo Fournier^(✉)  and Stanley Durrleman 

Sorbonne Université, Institut du Cerveau - Paris Brain Institute - ICM, CNRS, Inria,
Inserm, AP-HP, Hôpital de la Pitié Salpêtrière, 75013 Paris, France
`nemo.fournier@icm-institute.org`

Abstract. We introduce a disease progression model suited for neurodegenerative pathologies that allows to model associations between covariates and dynamic features of the disease course. We establish a statistical framework and implement an algorithm for its estimation. We show that the model is reliable and can provide uncertainty estimates of the discovered associations thanks to its Bayesian formulation. The model’s interest is showcased by shining a new light on genetic associations.

Keywords: Multimodal Disease Progression Modelling · Alzheimer’s Disease · Genetic Associations

1 Introduction

The clinical courses of neurodegenerative pathologies such as Alzheimer’s or Parkinson’s Diseases span multiple years and encompass intricate evolution of patients’ cognitive abilities, physiological biomarkers and brain structure. Longitudinal studies are an essential tool for clinicians to uncover the diseases’ mechanisms. In such studies, biomarkers and cognitive scores of patients are repeatedly measured at different times and need to be analyzed together, usually with a two-sided scope. First, to describe the general process at play across a whole cohort of patients: this is *population-level* modelling and allows to describe the average course of the disease. A second layer aims at explaining and predicting the variability observed among individuals: this is *personalized-level* modelling.

Mixed-effect frameworks are widely adopted to address these multi-layered prospects, offering to disentangle *fixed effects* (population level) from *random*

This work was funded in part by the French government under management of Agence Nationale de la Recherche as part of the Investissements d’avenir program, reference ANR-19-P3IA- 0001 (PRAIRIE 3IA Institute), ANR-19-JPW2-000 (E-DADS), and ANR10-IAIHU-06 (IHU ICM), as well as by the European Research council reference ERC-678304 and the H2020 programme via grant 826421 (TVB-Cloud).

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2023
H. Greenspan et al. (Eds.): MICCAI 2023, LNCS 14224, pp. 601–610, 2023.
https://doi.org/10.1007/978-3-031-43904-9_58

effects (individual level) to explain the variability of the disease. Linear mixed-effects models are the simplest instances of such models. Generalized linear and non-linear mixed-models are now often preferred to account for the neurodegenerative diseases' peculiarities, and most state of the art disease progression models (*e.g.* [2, 14, 16]) belong to these categories. They are indeed better suited to describe phenomena whose complex dynamic spans multiple years. They have been used with success to describe the natural history of diseases [8] or make individualized predictions, for instance to enrich clinical trials [11]. A general formulation is as follows, where η is a non linear mapping between timepoints and clinical markers, parametrized by fixed-effects α and random effects β_i (methods differ by the chosen non-linearity η and how α and β parametrize the disease course):

$$\mathbf{y}_i = \eta_{\alpha}(\mathbf{t}_i \mid \beta_i)$$

Inter-patient variability is thereby modelled through random perturbations β_i around a fixed reference α . However, it is known that some of this clinical variability between patients is explained by external factors (and thus hardly explained entirely by *random* perturbations). Genetic mutations or external factors such as gender, family history, education or socio-economics levels can influence the course of pathologies. Accumulating evidence suggests that the variability induced by such covariates stems from general mechanisms shared across the population, for instance in Alzheimer's and Parkinson's diseases [3, 6, 7, 10]. In the presented models, observed covariates \mathbf{c}_i are not taken into account and only the repeated observations \mathbf{y}_i are modelled as a function of the patient's ages \mathbf{t}_i . Thus, random effects might be such that $\mathbf{E}[\beta_i \mid \mathbf{c}_i] \neq 0$. This shows that some signal present in covariates to explain the progression of the disease has not been fully exploited.

Our contribution is to propose a slight change in this mixed-effect paradigm to allow non-linear models to also be influenced by these variables. Instead of estimating a fixed effect α (parametrizing the average disease course) as well as random effects β_i , we introduce a *link function* f_{φ} that can predict, given a set of covariates \mathbf{c}_i (*e.g.* sex, education level, SNP arrays, genetic risk scores), an expected trajectory of the disease conditioned by these covariates.

The main difference between the standard approach and our method is that the previously introduced fixed-effects α are now estimated for each subject as a deterministic function of their covariates $f_{\varphi}(\mathbf{c}_i)$. It also differs from accounting for the heterogeneity through hierarchical progression models [13, 17]) since covariates are, in our case, supervisingly used during model calibration and used to navigate through a *continuum* of disease models, instead of having defined clusters.

We demonstrate the value of this approach by adapting a general modelling framework, namely a non-linear Bayesian model: the *Disease Course Mapping* (DCM) [16]. We show that accounting for time-independent covariates in the longitudinal modelling with this approach can be done in a reasonable statistical setting. A stochastic estimation algorithm can be devised and we propose an instantiation and implementation of our model, which we validate first on

synthetic. We then use clinical data from the Alzheimer’s Disease Neuroimaging (ADNI) cohort and further demonstrate the clinical interest of the method by estimating new associations between genetics and disease dynamics.

2 Method

We derive here an algorithm that learns to model repeated observations while accounting for the heterogeneity explained by additional covariates.

2.1 A Generic Mixed-Effects Geometric Model

In their seminal paper [16], Schirrat *et al.* introduced a generic framework to model a dataset $(\mathbf{y}_{i,j})$ of multimodal longitudinal measurements. Here $\mathbf{y}_{i,j}$ is a vector of N biomarkers measured for the i -th subject at their j -th visit – *i.e.* at age $t_{i,j}$. Each observation $\mathbf{y}_{i,j}$ is assumed to lie on Riemannian submanifold (\mathcal{M}, g) of \mathbf{R}^N . The *average course* of the disease is posited to be such that individual progressions stem from a geodesic trajectory γ_0 on the manifold surface. Geodesic equations imply that γ_0 is entirely characterized by its initial position $p_0 \in \mathcal{M}$ and speed $v_0 \in T_{p_0}\mathcal{M}$ (tangent space of \mathcal{M} at p_0) at time t_0 . *Individual trajectories* are obtained from this reference trajectory γ_0 via a temporal reparametrization $t \mapsto \psi_i(t)$, used to derive what we name a *disease age* and enables registering patient’s chronological ages onto a common disease timeline. Spatial effects w_i are applied to the reference trajectory thanks to an exp-parallelization procedure that identifiably deforms geodesics in the manifold space. We denote $\eta_{\gamma_0}^{w_i}$ the resulting geodesic. We refer to [16] for extensive details on the geometric properties of these operations (such as commutativity and identifiability of both temporal and spatial effects).

The choice of the manifold’s metric shapes the geodesic trajectories and thus the disease model [4, 15, 16]. Clinical knowledge of Alzheimer’s Disease suggests that sigmoid shapes as sound candidates to model biomarkers’ evolutions (see [5] for clinical considerations, or [12] for the logistic dynamic of imaging-derived features such as brain-averaged protein loads backed by prion-like diffusion hypothesis). We therefore consider a product-metric g_p such that geodesic are sigmoids: $g_p(u, v) = u \cdot M(p) \cdot v$ with $M(p) = \frac{1}{p^2(1-p)^2}$, which gives the trajectories of Eq. (1). The resulting trajectories and the geometric interpretation of the (v_0, p_0, t_0) parameters are also presented in Fig. 1.

$$\eta_{\gamma_0}^{w_i}(\psi_i(t_{i,j}))^{(k)} = \left(1 + \left(\frac{1}{p_0^{(k)}} - 1 \right) \exp \left(- \frac{v_0^{(k)} \psi_i(t_{i,j}) + w_i^{(k)}}{p_0^{(k)} (1 - p_0^{(k)})} \right) \right)^{-1} \quad (1)$$

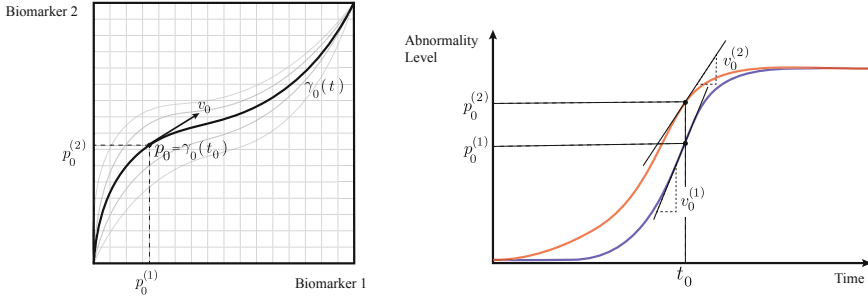


Fig. 1. A two feature model, with geodesic trajectories on the manifold (left) and the biomarkers observation space (right). This provides the intuition over the effect of the initial position p_0 and the initial velocity v_0 at time t_0 .

2.2 Covariate Association and Statistical Framework

We provide here a statistical instantiation of the previous geometric model. As described, given a geodesic trajectory γ_0 (fully specified by its position p_0 and speed v_0 at initial time t_0), and a set of random effect ψ_i and w_i , individual trajectories of an individual i observed at times $(t_{i,j})_j$ are modelled by the curve $\eta_{\gamma_0}^{w_i}(\psi_i(t_{i,j}))$. We propose that γ_0 (which represents the reference disease course, as a fixed-effect of the model) is to be computed for each subject i from the measured covariates c_i as:

$$\gamma_{0,(i)} \cong (p_{0,(i)}, v_{0,(i)}, t_{0,(i)}) = f_{\varphi}(c_i)$$

where f belongs to a parametrized family of functions and φ are its parameters treated as the new fixed-effect of the model. The individual effects to register this computed $\gamma_{0,(i)}$ onto observations are characterized by two random effects: an acceleration factor ξ_i and a time-shift τ_i such that $\psi_i(t) = e^{\xi_i} (t - t_{0,(i)} - \tau_i)$. On top of these are space-shifts $w_i \in \mathbf{R}^N$, computed thanks to an ICA: $w_i = A s_i$, where A is a latent matrix of independent directions (fixed effect) and s_i is the corresponding individual latent source vector (random effect).

Our hierarchical statistical model treats the fixed and random effects as a set of latent variables z which is the reunion of the population and individual variables $z_{\text{pop}} = \{\varphi, A\}$ and $z_{\text{indiv}} = \{(s_i)_i, (\tau_i)_i, (\xi_i)_i\}$. We posit the following priors on these latent parameters, where $\theta_{\text{hyper}} = \{\sigma_{\varphi}, \sigma_A\}$ are fixed hyperparameters and $\theta_{\text{model}} = \{\bar{\varphi}, \bar{A}, \sigma_{\tau}^2, \sigma_{\xi}^2, \sigma^2\}$ are the parameters of the model to be estimated:

$$\varphi \sim \mathcal{N}(\bar{\varphi}, \sigma_{\varphi}^2) \quad A \sim \mathcal{N}(\bar{A}, \sigma_A^2) \quad \xi_i \sim \mathcal{N}(0, \sigma_{\xi}^2) \quad \tau_i \sim \mathcal{N}(0, \sigma_{\tau}^2) \quad s_i \sim \mathcal{N}(0, 1)$$

A non-informative prior is used over these model parameters due to the lack of a-priori knowledge. We seek to maximize a posteriori the joint-likelihood under the following additive Gaussian noise modelling $y_{i,j} = \eta_{f_{\varphi}(c_i)}^{A s_i}(\psi_i(t_{i,j})) + \varepsilon_{i,j}$

$$q(\mathbf{y}, \mathbf{z}, \theta_{\text{model}}) = \underbrace{q(\mathbf{y} \mid \mathbf{z}, \theta_{\text{model}})}_{\text{data attachment}} \times \underbrace{q(\mathbf{z} \mid \theta_{\text{model}}, \theta_{\text{hyper}})}_{\text{regularization of the latent variables}} \times \underbrace{q_{\text{prior}}(\theta_{\text{model}} \mid \theta_{\text{hyper}})}_{\text{prior on model parameters (taken non informative)}}$$

It can be shown that the model's likelihood function lies in the curved exponential family. That is there exist two smooth functions Φ and Ψ functions of θ_{model} and a measurable sufficient statistics function $S(\mathbf{y}, \mathbf{z})$ of the data and the latent realizations such that $\log q(\mathbf{y}, \mathbf{z}, \theta_{\text{model}})$ factors as:

$$\log q(\mathbf{y}, \mathbf{z}, \theta_{\text{model}}) = -\Phi(\theta_{\text{model}}) + \langle S(\mathbf{y}, \mathbf{z}), \Psi(\theta_{\text{model}}) \rangle$$

This allows estimating our model with a Monte-Carlo Markov-Chain Stochastic Approximation version of the Expectation Maximization algorithm (MCMC-SAEM) while enjoying theoretical guarantees of convergence [1]. The expectation phase is therefore built upon a sampling scheme to sample from the posterior distribution of the latent parameters (namely Metropolis-Hastings within Gibbs sampler). The maximization phase follows update rules established by finding critical points of $\theta \mapsto -\Phi(\theta) + \langle \tilde{S}_p, \Psi(\theta) \rangle$ (\tilde{S}_p is the stochastic approximation of the sufficient statistics built at step p), which yields analytic expressions.

We choose to parametrize the link function as a linear mapping between covariates and dynamic parameters $f_{\varphi}(\mathbf{c}_i) = \varphi_{\text{slope}} \cdot \mathbf{c}_i + \varphi_{\text{intercept}}$. This will provide an interpretable model to explain the general processes linking covariates to dynamic features such as the base pace of the disease or average onset time. The coefficients of φ that correspond to the mapping between the covariates \mathbf{c}_i and v_0 measure how much a given covariate impact the progression speed of each feature, and can be analyzed easily. Model parameters are initialized by setting their intercept to the models learned by a regular DCM model without considering covariates, while latent parameters are initialized at random.

3 Evaluation, Clinical Results and Discussion

3.1 Simulated Data

We used the generative abilities of the DCM[8] to simulate multimodal longitudinal datasets with covariates influencing the dynamic of the progression. To this end, we fixed some reference models corresponding each to a slightly different *pure form* of a fictional disease. Then, covariates were simulated either:

- as binary covariates that directly dictated which hardcoded model is used to simulate the repeated measurements (covariate thought as a mutation-status or sex for instance).
- as continuous covariates, influencing the simulated progression by using them as convex coefficients in a combination of the reference models. The covariates are seen as continuous risk factors of following one form or another.

Our simulated datasets typically included 500 subjects with an average of 5 visits and an average follow-up duration of about 5 ± 2 years and a measurement noise of around 5%. Such an experiment is summarized in Fig. 2, where three continuous covariates were simulated on the $[0, 1]$ range, the first one being a risk to develop a motor form of the disease, the second one a memory-form risk and the third a covariate without any influence on the disease.

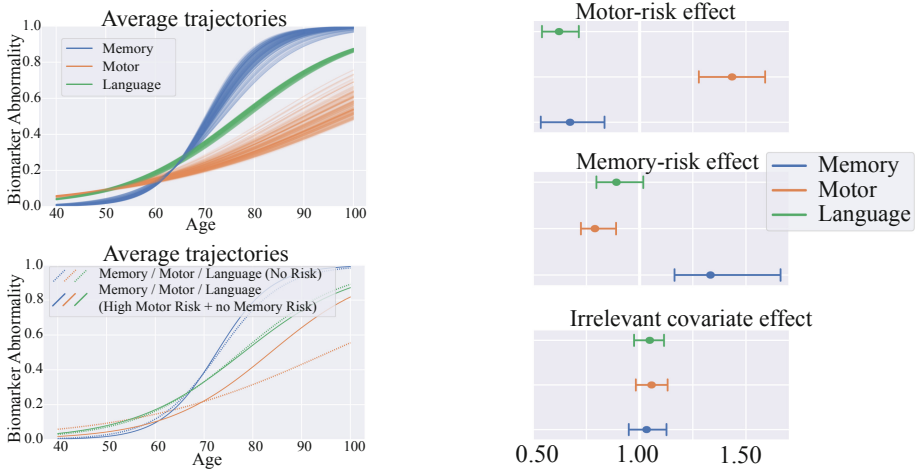
In this example, our calibrated model correctly matches each covariate to its simulated effect. For instance: the coefficients of the link function related to the disease initial speed v_0 (we refer to Fig. 1 for its interpretation) associated with the memory-risk covariate show an acceleration of the decline in memory (multiplicative factor of 1.32 [1.15, 1.62] - credible interval at 95%) contrasted to the two other biomarkers (factor of 0.79 [0.73, 0.90] for the motor and 0.89 [0.80, 1.01] for the language). These intervals are represented in Fig. 2b. These two other features are slowed down relatively to the memory in order for the model to capture the change in the slope-ratio of different features on the fixed effects. If we translate the effects of these coefficients into effects on *slope-ratio*, we obtain indeed obtain that the ratio between memory-speed and motor-speed goes from 4.22 [3.54, 5.02] (no extra memory-risk) to 7.05 [5.54, 11.26] (maximum extra memory-risk). The ground truth change of slope (from the reference models) was from 4.48 (no particular risk) to 9.19 (full risk) and is therefore covered by our credible intervals.

Similarly, the coefficients associated with the irrelevant covariate did not capture any significant effect (factors of 1.03 [0.94, 1.12], 1.05 [0.98, 1.12] and 1.04 [0.97, 1.11] for the memory, motor and language features), which validates the ability of the model to discard covariates without influence on the disease dynamic.

3.2 Multimodal Clinical Data

Data used in the preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu). ADNI was launched in 2003 as a public-private partnership led by Michael W. Weiner, MD. For up-to-date information, see <https://www.adni-info.org>. We selected subjects that eventually converted to an MCI or AD stage during their follow-up. This amounted to 1440 patients for a total of 9343 visits. The follow-up duration was 4.069 (± 3.190) years, with a baseline age of 73.683 (± 7.508) years old. We processed and included biomarkers relevant to monitor AD progression:

- Two cognitive scores: the Mini-Mental State Exam (MMSE) and the AD Assessment Scale-Cognitive (ADAS-Cog). We normalized and inverted them so that they both cover the $[0, 1]$ interval, (1 being the highest abnormality).
- Hippocampus and Ventricles volumes, measured by structural T1 MRI and normalized by patient's Intracranial Volume (ICV). As for the cognitive scores, these measurements were rescaled to $[0, 1]$ interval.
- Contrasted PET imaging derived brain-averaged amyloid $\beta 42$ and phosphorylated τ proteins loads, also rescaled to $[0, 1]$.



(a) *Top.* Continuum of models to simulate data. *Bottom.* recovered disease course and effect of the motor-risk covariate.

(b) Posterior mode and credible intervals on the *covariates - progression - speed* interaction coefficient

Fig. 2. Risk factors to develop an acute memory-form or acute motor-form of the disease are sampled continuously from the $[0, 1]$ interval, and used as convex coefficients to combine three reference models (a standard form, a motor-dominant form and an acute memory form), yielding the continuum of possible trajectories presented in (a). We also sample an *irrelevant* covariate that is never used to modulate the disease course. In (a) is the resulting model, which we can visualize for any combination of covariates (two combinations are presented). We also plot the credible intervals (at 95%) for the coefficients linking covariates to the speed of progression on each feature (multiplicative effect, thus 1.0 stands for no influence while a coefficient of 1.5 stands for an expected progression speed greater by 50%).

APOE- $\epsilon 4$. We calibrated our model by including a covariate known to modulate Alzheimer’s Disease course, namely the patient’s APOE mutation status. The results of this model are showcased in Fig. 3. It shows that the APOE mutation, which is a known risk factor for AD, has a clear effect on the disease dynamic: in the obtained disease course map, the mutation is associated with earlier and faster abnormalities on most biomarkers. We also investigate the learned linked function f_φ and its coefficients dictating the interaction between the covariates and the speed of progression. This is presented in Fig. 3. This showcases how the contribution of the APOE to the speed of progression (as in the coefficient linking the covariate to the coordinates of v_0) is different among biomarkers. Features can be grouped into a *cognitive scores* group, more impacted than the other features by the mutation (1.49 [1.31, 1.66] and 1.34 [1.21, 1.46] respectively for MMSE and ADAS), a structural subgroup that also shows significative (even though slower) increase of progression speed (1.08 [1.00, 1.14] and 1.07 [1.0, 1.15] resp. for Hippocampus and Ventricles volume), while exhibiting less clear effects

on the proteins loads (1.01 [0.91, 1.12] and 1.03 [0.92, 1.18] resp for Amyloid and tau loads).

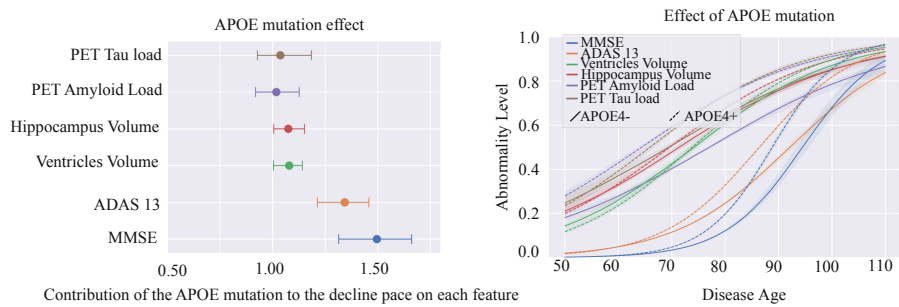


Fig. 3. Model of Alzheimer’s Disease obtained from ADNI data. Left: estimated parameters of the link function φ (mode of the posterior and 95% credible interval). Right: difference in the trajectory conditioned by the mutation status (the represented trajectories are for 0 copy vs 2 copies of the APOE- ϵ 4 allele).

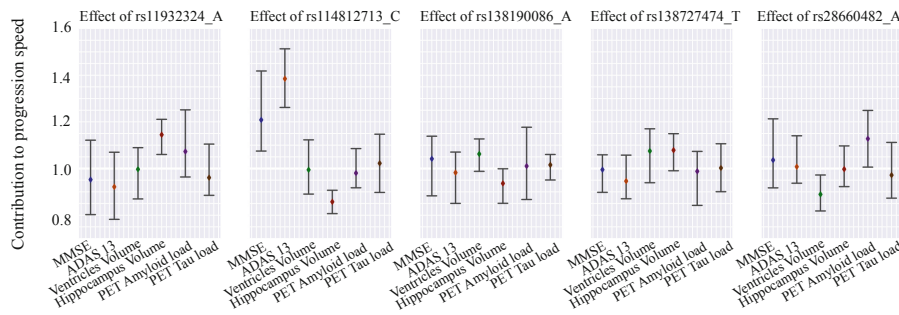


Fig. 4. Analysis of the learned interaction between some of the SNP included in the analysis and the speed of progression of each of the 6 measured features. Some SNPs present no significant interaction with the progression speed of any of the variables (*e.g.* rs138727474T) to SNPs that are associated with a group of feature (*e.g.* cognitive domain for rs114812713C) or single features (*e.g.* rs11932324A or rs28660482A, in either a protective or risk-inducing direction).

SNP Associations. We selected a subset of 69 Single Nucleotide Polymorphisms (SNP) among the top associations with AD diagnosis from a reference Genome-Wide Association Study (GWAS) [9]. We included them in our model as covariates. The results suggest that being associated with the diagnosis does not inform a priori on the influence of each SNP on the disease course. In Fig. 4 we show

that, even though all these SNP were selected for being significantly associated with the diagnosis, they can exhibit differences in their association with the disease dynamic.

4 Conclusion

We proposed a framework to adapt a state of the art Bayesian non-linear mixed-effect disease progression model to capture the effects of external covariates into the disease dynamic. We implemented an estimation algorithm, and show that it reliably provides new interpretable measures of interaction between covariates and the disease course. For instance, we recover the (clinically known) association between the APOE- ϵ 4 mutation and cognitive dysfunction. In particular, its use on genetic data (either single mutation status or SNP arrays) could help to go beyond associations with the sole diagnosis and provide complementary tools to GWAS.

References

1. Allasonnière, S., Kuhn, E., Trouvé, A.: Construction of Bayesian deformable models via a stochastic approximation algorithm: a convergence study. *Bernoulli* **16**(3), 641–678 (2010). <https://doi.org/10.3150/09-BEJ229>
2. Donohue, M.C., et al.: Estimating long-term multivariate progression from short-term data. *Alzheimer's & Dementia* **10**(5S), S400–S410 (2014). <https://doi.org/10.1016/j.jalz.2013.10.003>
3. Greenland, J.C., Williams-Gray, C.H., Barker, R.A.: The clinical heterogeneity of Parkinson's disease and its therapeutic implications. *Eur. J. Neurosci.* **49**(3), 328–338 (2019). <https://doi.org/10.1111/ejn.14094>
4. Gruffaz, S., Poulet, P.E., Maheux, E., Jedynek, B., Durrleman, S.: Learning Riemannian metric for disease progression modeling. In: *Advances in Neural Information Processing Systems*, vol. 34, pp. 23780–23792. Curran Associates, Inc. (2021)
5. Jack, C.R., et al.: Tracking pathophysiological processes in Alzheimer's disease: an updated hypothetical model of dynamic biomarkers. *Lancet Neurol.* **12**(2), 207–216 (2013). [https://doi.org/10.1016/S1474-4422\(12\)70291-0](https://doi.org/10.1016/S1474-4422(12)70291-0)
6. Jutten, R.J., Sikkes, S.A., Van der Flier, W.M., Scheltens, P., Visser, P.J., Tijms, B.M.: for the Alzheimer's disease neuroimaging initiative: finding treatment effects in Alzheimer trials in the face of disease progression heterogeneity. *Neurology* **96**(22), e2673–e2684 (2021). <https://doi.org/10.1212/WNL.00000000000012022>
7. Komarova, N.L., Thalhauser, C.J.: High degree of heterogeneity in Alzheimer's disease progression patterns. *PLoS Comput. Biol.* **7**(11), e1002251 (2011). <https://doi.org/10.1371/journal.pcbi.1002251>
8. Koval, I., et al.: AD course map charts Alzheimer's disease progression. *Sci. Rep.* **11**(1), 8020 (2021). <https://doi.org/10.1038/s41598-021-87434-1>
9. Kunkle, B.W., et al.: Genetic meta-analysis of diagnosed Alzheimer's disease identifies new risk loci and implicates A β , tau, immunity and lipid processing. *Nat. Genet.* **51**(3), 414–430 (2019). <https://doi.org/10.1038/s41588-019-0358-2>
10. Livingston, G., et al.: Dementia prevention, intervention, and care: 2020 report of the lancet commission. *The Lancet* **396**(10248), 413–446 (2020). [https://doi.org/10.1016/S0140-6736\(20\)30367-6](https://doi.org/10.1016/S0140-6736(20)30367-6)

11. Maheux, E., et al.: Forecasting individual progression trajectories in Alzheimer's disease. *Nat. Commun.* **14**(1), 761 (2023). <https://doi.org/10.1038/s41467-022-35712-5>
12. Meisl, G., et al.: In vivo rate-determining steps of tau seed accumulation in Alzheimer's disease. *Sci. Adv.* **7**(44), eabh1448 (2021). <https://doi.org/10.1126/sciadv.abh1448>
13. Poulet, P.-E., Durrleman, S.: Mixture modeling for identifying subtypes in disease course mapping. In: Feragen, A., Sommer, S., Schnabel, J., Nielsen, M. (eds.) *IPMI 2021. LNCS*, vol. 12729, pp. 571–582. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-78191-0_44
14. Raket, L.L.: Statistical disease progression modeling in alzheimer disease. *Front. Big Data* **3** (2020). <https://doi.org/10.3389/fdata.2020.00024>
15. Sauty, B., Durrleman, S.: Riemannian metric learning for progression modeling of longitudinal datasets. In: 2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI), pp. 1–5 (Mar 2022). <https://doi.org/10.1109/ISBI52829.2022.9761641>
16. Schiratti, J.B., Allasonnière, S., Colliot, O., Durrleman, S.: A Bayesian mixed-effects model to learn trajectories of changes from repeated manifold-valued observations. *J. Mach. Learn. Res.* **18**(1), 4840–4872 (2017)
17. Young, A.L., et al.: Uncovering the heterogeneity and temporal complexity of neurodegenerative diseases with subtype and stage inference. *Nat. Commun.* **9**(1), 4273 (2018). <https://doi.org/10.1038/s41467-018-05892-0>