Interpretable discriminant analysis for functional data supported on random nonlinear domains with an application to Alzheimer's disease

Eardi Lila<sup>1</sup>, Wenbo Zhang<sup>1,2</sup>, Swati Rane Levendovszky<sup>3</sup>, for the Alzheimer's Disease Neuroimaging Initiative\*

<sup>1</sup>Department of Biostatistics, University of Washington, USA <sup>2</sup>Department of Statistics, University of California, Irvine, USA <sup>3</sup>Department of Radiology, University of Washington, USA

#### Abstract

We introduce a novel framework for the classification of functional data supported on nonlinear, and possibly random, manifold domains. The motivating application is the identification of subjects with Alzheimer's disease from their cortical surface geometry and associated cortical thickness map. The proposed model is based upon a reformulation of the classification problem as a regularized multivariate functional linear regression model. This allows us to adopt a direct approach to the estimation of the most discriminant direction while controlling for its complexity with appropriate differential regularization. Our approach does not require prior estimation of the covariance structure of the functional predictors, which is computationally prohibitive in our application setting. We provide a theoretical analysis of the out-of-sample prediction error of the proposed model and explore the finite sample performance in a simulation setting. We apply the proposed method to a pooled dataset from the Alzheimer's Disease Neuroimaging Initiative and the Parkinson's Progression Markers Initiative. Through this application, we identify discriminant directions that capture both cortical geometric and thickness predictive features of Alzheimer's disease that are consistent with the existing neuroscience literature.

#### 1 Introduction

Functional discriminant analysis, a statistical framework used to predict categorical outcomes from functional predictors, has been extensively studied within the field of Functional Data Analysis (FDA) (Ramsay and Silverman 2015; Ferraty and Vieu 2006; Horváth and Kokoszka 2012; Hsing and Eubank 2013) and has motivated a large body of literature (see, e.g., James and Hastie 2001; Müller 2005; Preda 2007; Delaigle and Hall 2012; Dai et al. 2017; Berrendero et al. 2018; Kraus and Stefanucci 2019; Park et al. 2021). However, most of the existing methods are concerned with the classification of functions supported on one-dimensional linear domains, which can be a limiting assumption in many modern biomedical applications (Zhu et al. 2023). On the other hand, recent work on the analysis of functional data with manifold structure has mostly focused on vector-valued functions with non-Euclidean constraints in the image space (see e.g., Su et al. 2014; Dai and Müller 2018; Lin et al. 2017; Dubey and Müller 2020; Kim et al. 2021).

In this paper, motivated by the analysis of modern multi-modal imaging data, we propose a novel functional discriminant analysis model that can handle functional predictors supported on nonlinear

<sup>\*</sup>Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (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

sample-specific manifold domains, which we term Functions on Surfaces (FoSs) (Lila and Aston 2020). An example of such 'object data' (Marron and Dryden 2021) is provided in Figure 1A, illustrating our motivating application of identifying subjects with Alzheimer's disease from FoSs that are subject-specific cortical surfaces coupled with cortical thickness measurements. The statistical analysis of these object data poses unique statistical challenges. This is mainly due to the non-Euclidean structure of each individual domain, which makes it difficult to define appropriate spatial regularization, and due to the more abstract non-Euclidean structure of the latent space where the random domains are supported, which further invalidates traditional linear statistical models (Kendall 1984; Grenander and Miller 1998; Dryden and Mardia 2016; Younes 2019).

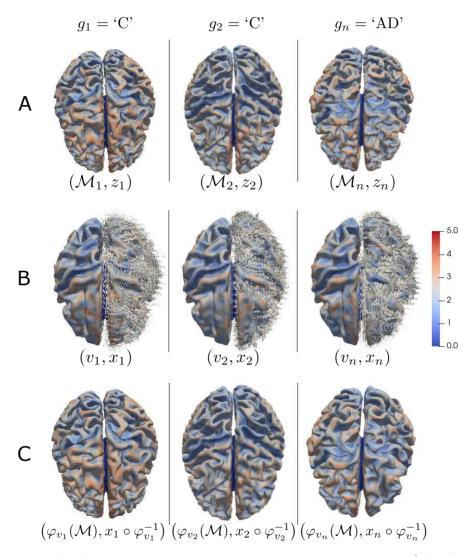


Figure 1: Panel A: FoSs of three subjects in the training sample, where  $g_i \in \{\text{`C'}, \text{`AD'}\}$  denotes the disease state of the *i*th individual (C: Control, AD: Alzheimer's Disease),  $\mathcal{M}_i$  is a two-dimensional manifold encoding the geometry of the cerebral cortex, and  $z_i : \mathcal{M}_i \to \mathbb{R}$  is a real function, supported on  $\mathcal{M}_i$ , describing cortical thickness (in mm). Panel B: Linear representation  $(v_i, x_i)$  of each FoS  $(\mathcal{M}_i, z_i)$  shown in Panel A. Here  $v_i : \mathbb{R}^3 \to \mathbb{R}^3$  is a vector-valued function encoding the geometry of the *i*th individual. This is depicted as a collection of 3D vectors  $\{v_i(p_j)\}$  for a dense set of points  $\{p_j\} \subset \mathbb{R}^3$ . For clarity, the function  $v_i$  is displayed only on half of its domain  $\mathbb{R}^3$ . The function  $x_i : \mathcal{M} \to \mathbb{R}$  describes the spatially normalized thickness map of the *i*th individual on the fixed template  $\mathcal{M}$ . Panel C: FoS  $(\varphi_{v_i}(\mathcal{M}), x_i \circ \varphi_{v_i}^{-1})$  parametrized by the associated functions  $(v_i, x_i)$  in Panel B. This is a close approximation of the FoS  $(\mathcal{M}_i, z_i)$  in Panel A.

Current approaches in the literature do not comprehensively address these challenges. Although

various methods have been proposed to model functional data on multi-dimensional domains, they often focus on flat domains (Goldsmith et al. 2014; Wang and Zhu 2017; Kang et al. 2018; Feng et al. 2020; Yu et al. 2021) or assume nonlinear but fixed domains (Chung et al. 2015; Chung et al. 2021; Lila et al. 2016; Mejia et al. 2020).

There has also been considerable work on the simpler setting of random surfaces that are *not* coupled with functional data. These efforts can be broadly grouped into three main approaches. The first approach leverages global parametrizations to represent surfaces, employing either an  $L^2$  metric (Chung et al. 2008; Epifanio and Ventura-Campos 2014; Ferrando et al. 2020) or a non-Euclidean metric (Jermyn et al. 2012; Jermyn et al. 2017; Kurtek and Drira 2015; Zhang et al. 2022). The second approach, which is more closely related to the one adopted in this work, uses diffeomorphic deformation functions of the surfaces' embedding space (Vaillant et al. 2004; Younes 2019; Arguillère et al. 2016), allowing for the inclusion of topological constraints. The third approach, prevalently used in neuroimaging studies, employs pre-specified or spectrum-based descriptors of shape (Reuter et al. 2006; Im et al. 2008; Wachinger et al. 2015; Hazlett et al. 2017; Wang and Wang 2017; Dong et al. 2019; Dong et al. 2020). A critical drawback of the latter approach is the inability to uniquely map the discrete representations back to the original space of random surfaces.

The statistical analysis of random surfaces that are *coupled* with functional data has not been extensively explored. One exception is the model in Zaetz and Kurtek (2015) which focuses on annotated surfaces. In addition, unsupervised models have been investigated by Charlier et al. (2017) and Lila and Aston (2020). Lee et al. (2017) have dealt with the classification problem by employing the fshape framework (Charlier et al. 2017) to represent the data and by using a linear discriminant model on the resulting finite-dimensional representations. In contrast, the statistical framework for the discriminant analysis of FoSs presented in this paper avoids any dimension reduction of the predictors and instead employs spatial penalties to regularize the discriminant direction. It tackles the non-Euclidean nature of the latent space of random domains by defining appropriate linear functional representations of FoSs, effectively reframing the problem of classifying FoSs as the problem of classifying bivariate functional data supported on general, but fixed, domains. A direct model for the estimation of a functional discriminant direction is then defined on the representation space where differential spatial regularizations are introduced to produce interpretable and well-defined estimates. A key feature of the proposed representation is its invertibility, which enables us to map estimates from the representation space back to the original space. This allows us to explore and interpret the estimated classification rule in the context of the original neurobiological objects of our motivating application.

The rest of the paper is organized as follows. In Section 2, we describe the representation model adopted to parametrize the non-Euclidean space of FoSs using linear function spaces. In Sections 3 and 4, we develop a novel discriminant analysis model on the parametrizing linear function spaces and provide theoretical guarantees for the prediction performance of the proposed model. We introduce an efficient discretization approach in Section 5 and apply the proposed model to the Alzheimer's Disease Neuroimaging Initiative (ADNI) and Parkinson's Progression Markers Initiative (PPMI) datasets in Section 6. Proofs and simulations are left to the appendices.

### 2 Functional data supported on general random domains

The data considered in this work is a sample of triplets

$$\{(g_i, \mathcal{M}_i, z_i), i = 1, \dots, n\},$$
 (1)

where  $g_i$  is a binary label,  $\mathcal{M}_i \subset \mathbb{R}^3$  is a sample-specific closed two-dimensional manifold embedded in  $\mathbb{R}^3$ , and  $z_i : \mathcal{M}_i \to \mathbb{R}$  is a scalar function supported on the geometric object  $\mathcal{M}_i$ . We moreover assume that the points on the geometries  $\{\mathcal{M}_i\}$  of the observed FoSs are in one-to-one correspondence across subjects. We refer to the pairs  $\{(\mathcal{M}_i, z_i)\}$  as FoSs.

In Figure 1A, we display three observations of the training sample of our final application, where  $g_i$  encodes the disease state of the subject,  $\mathcal{M}_i$  encodes the geometry of the cerebral cortex, and  $z_i$  encodes the cortical thickness map supported on  $\mathcal{M}_i$ . Our goal is to build a classifier from the given training sample that can predict the binary label  $g^*$  of a previously unseen FoS  $(\mathcal{M}^*, z^*)$ .

#### 2.1 Linear functional representations

In our motivating application, the geometric objects  $\{\mathcal{M}_i : i = 1, ..., n\}$  are topologically equivalent to a sphere, and therefore, they do not display holes or self-intersections. To inform our model of such physical non-Euclidean constraints, we define a convenient unconstrained representation model for the FoSs  $\{(\mathcal{M}_i, z_i) : i = 1, ..., n\}$  in terms of objects belonging to linear function spaces.

Let  $\mathcal{M}$  be a template two-dimensional manifold embedded in  $\mathbb{R}^3$  that is topologically equivalent to a sphere. We denote by  $\mathcal{L}^2(\mathcal{M})$  the space of square integrable functions over  $\mathcal{M}$ , equipped with the standard inner product  $\langle \cdot, \cdot \rangle_{\mathcal{L}^2(\mathcal{M})}$  and norm  $\|\cdot\|_{\mathcal{L}^2(\mathcal{M})}$ , and denote by  $\mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3)$  the space of square integrable vector-valued functions from  $\mathbb{R}^3$  to  $\mathbb{R}^3$ , with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3)}$  and norm  $\|\cdot\|_{\mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3)}$ . Let  $\mathcal{V}(\mathbb{R}^3) \subset \mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3)$  be a Reproducing Kernel Hilbert Space (RKHS) of smooth functions with compact support in  $\mathbb{R}^3$ . We then introduce a diffeomorphic operator  $\varphi$  such that  $\varphi_v : \mathbb{R}^3 \to \mathbb{R}^3$  is a diffeomorphic function for every choice of  $v \in \mathcal{V}(\mathbb{R}^3)$  (Younes 2019). We denote by  $\varphi_v(\mathcal{M})$  the geometric object given by displacing every point  $p \in \mathcal{M} \subset \mathbb{R}^3$  to the new location  $\varphi_v(p)$ . A direct consequence of  $\varphi_v$  being diffeomorphic is that  $\varphi_v(\mathcal{M})$  is topologically equivalent to a sphere for every choice of  $v \in \mathcal{V}(\mathbb{R}^3)$ . The construction of the diffeomorphic operator adopted in this paper and the computation of  $v_i \in \mathcal{V}(\mathbb{R}^3)$  such that  $\varphi_{v_i}(\mathcal{M}) \approx \mathcal{M}_i$  are detailed in Appendix A.1.

Next, we use the estimated  $v_i$ , and  $z_i$ , to define the spatially normalized function  $x_i : \mathcal{M} \to \mathbb{R}$  as  $x_i = z_i \circ \varphi_{v_i}$ . For each point  $p \in \mathcal{M}$ ,  $x_i(p)$  is given by  $z_i \circ \varphi_{v_i}(p)$ , that is, the value of  $z_i$  at the corresponding point  $\varphi_{v_i}(p) \in \mathcal{M}_i$ . We can then represent each FoS  $(\mathcal{M}_i, z_i)$  with a pair of functions  $(v_i, x_i)$  such that

$$(\mathcal{M}_i, z_i) \approx (\varphi_{v_i}(\mathcal{M}), x_i \circ \varphi_{v_i}^{-1}),$$
 (2)

where  $v_i \in \mathcal{V}(\mathbb{R}^3)$  and  $x_i \in \mathcal{L}^2(\mathcal{M})$ .

We can depict the representation model introduced as follows:

$$(v_i, x_i) \longleftrightarrow (\mathcal{M}_i, z_i),$$
 (3)

meaning that given a FoS  $(\mathcal{M}_i, z_i)$ , we can compute a loss-less representation  $(v_i, x_i)$  as described earlier, and vice-versa, given the representation  $(v_i, x_i)$ , we can compute the associated FoS through equation (2). Hence, the pair of functions  $(v_i, x_i)$  provides us with a linear representation of the original FoS  $(\mathcal{M}_i, z_i)$  where every geometric object  $\mathcal{M}_i$  is modeled as a (diffeomorphic) deformation of the template, i.e.  $\varphi_{v_i}(\mathcal{M})$ , while the associated function  $z_i$  is given by 'transporting' the spatially normalized function  $x_i$  onto  $\mathcal{M}_i$  with such a deformation.

The approach described allows us to recast the original non-Euclidean problem of learning a classifier from the training sample  $\{(g_i, \mathcal{M}_i, z_i) | i = 1, ..., n\}$  as the problem of learning a classifier from

$$\{(q_i, v_i, x_i)|i=1,\dots,n\},$$
 (4)

where  $v_i \in \mathcal{V}(\mathbb{R}^3) \subset \mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3)$  and  $x_i \in \mathcal{L}^2(\mathcal{M})$  are two functional predictors both belonging to *linear* function spaces. In Figure 1B, we show the functional linear representations associated with the three FoSs in Figure 1A.

Crucially, the representation mapping employed here is 'invertible', meaning that any pair of estimates  $(\beta^G, \beta^F) \in \mathcal{V}(\mathbb{R}^3) \times \mathcal{L}^2(\mathcal{M})$ , such as the 'direction' that optimally discriminates between two classes, can be mapped back to the original space of FoS using equation (2). This mapping defines the associated trajectories of FoSs

$$\left\{ \left( \varphi_{c_1\beta^G}(\mathcal{M}), c_2\beta^F \circ \varphi_{c_1\beta^G}^{-1} \right), c_1, c_2 \in \mathbb{R} \right\}, \tag{5}$$

where  $\varphi_{c_1\beta G}(\mathcal{M})$  is guaranteed to be topologically equivalent to a sphere, thereby satisfying the physical constraints of the problem considered.

In contrast to methods that require computing shape features, such as the spectrum of the Laplace-Beltrami operator (Reuter et al. 2006; Wachinger et al. 2015), the approach adopted here provides us with interpretable discriminant directions in the space of the original neurobiological objects, as described by equation (5). In addition, unlike approaches that work with global parametrizations of FoSs (see, e.g., Chung et al. 2008; Epifanio and Ventura-Campos 2014; Zaetz and Kurtek 2015; Ferrando et al. 2020), the representation model used in this study is independent of the imaging data type, as long as we can specify how  $\varphi_{v_i}$  deforms our objects and a suitable similarity measure. Therefore, the framework proposed in this work has the potential to accommodate additional types of data, such as streamlines generated from diffusion tensor images, where there may not be a one-to-one correspondence across subjects, but for which optimal transport similarity measures have been developed (Feydy et al. 2017). Although our work assumes that the FoSs are in one-to-one correspondence, this is not strictly necessary. Assuming a one-to-one correspondence simplifies the definition of a similarity measure and makes it easier to compute these representations for complex objects such as cortical surfaces, as detailed in Appendix A.1. In Section 6, we compare the performance of our representation model with alternative models, in the context of our motivating application.

#### 2.2 Discriminant analysis on the parametrizing linear function spaces

The aim of Sections 3 and 4 is to provide methodology for learning a linear classifier, from the training data  $\{(g_i, v_i, x_i) | i = 1, ..., n\}$  displayed in Figure 1B, by introducing a novel functional classification model that has the following crucial characteristics:

- Does not rely on Functional Principal Components Analysis (FPCA), or related dimension reduction methods, to reduce the dimension of the functional predictors, bypassing the intrinsic assumption that the discriminant direction is well represented by the space spanned by the first few unsupervised PC functions;
- Can be applied to bivariate, and possibly multivariate, functional predictors each supported on a different domain;
- Allows for explicit spatial regularization of the estimates on potentially nonlinear manifold domains, yielding well-defined and interpretable estimates;
- Provides a direct approach to estimating the discriminant directions without relying on prior computation of the covariance structure, which is prohibitive in our application setting.

## 3 Linear discriminant analysis over general domains

We begin by focusing on the sub-problem of defining a classifier for the training data  $\{(g_i, x_i)\}$ , i.e., for the spatially normalized functional predictors supported on the fixed nonlinear manifold  $\mathcal{M}$ . In Section 4, we then extend the proposed model to account for the geometric component  $v_i$ , in an additive fashion.

Assume the training set  $\{(g_i, x_i)\}$  consists of n independent copies of (G, X), a pair of random variables with X a zero-mean random function taking values in  $\mathcal{L}^2(\mathcal{M})$  and G a binary random variable such that  $P(G=1)=\pi_1$  and  $P(G=2)=\pi_2$ . Let  $\mu_1=\mathbb{E}\left[X|G=1\right]$  and  $\mu_2=\mathbb{E}\left[X|G=2\right]$  denote the conditional means of X and assume  $\mu_1\neq\mu_2$ . Moreover, let  $C(p,q)=\mathbb{E}\left[X(p)X(q)\right]$ ,  $p,q\in\mathcal{M}$  denote the covariance function of X and assume this is square integrable, i.e.,  $\int_{\mathcal{M}}\int_{\mathcal{M}}C(p,q)^2\,dp\,dq<\infty$ . For a real, symmetric, square-integrable, and non-negative function  $R\in\mathcal{L}^2(\mathcal{M})\to\mathcal{L}^2(\mathcal{M})$  be defined as

$$L_R(\beta)(\cdot) = \int_{\mathcal{M}} R(p, \cdot)\beta(p) dp, \, \forall \beta \in \mathcal{L}^2(\mathcal{M}).$$

Consequently,  $L_C$  denotes the covariance operator of X, which is a compact self-adjoint operator and therefore admits the following spectral representation

$$L_C(\beta) = \sum_{k=1}^{\infty} \theta_k \langle \beta, e_k \rangle_{\mathcal{L}^2(\mathcal{M})} e_k, \tag{6}$$

in terms of the eigenvalues  $\theta_1 \geq \theta_2 \geq \ldots \geq 0$  and associated eigenfunctions  $e_1, e_2, \ldots \subset \mathcal{L}^2(\mathcal{M})$  of  $L_C$ . Let  $L_C^{-1}$  denote the linear inverse covariance operator, where  $L_C^{-1}(e_k) = \frac{1}{\theta_k} e_k$  for all  $k \geq 1$ . Assume that  $\|L_C^{-1}(\mu_2 - \mu_1)\|_{\mathcal{L}^2(\mathcal{M})} < \infty$  and define the population quantity  $\beta^0 \in \mathcal{L}^2(\mathcal{M})$  such that

$$L_C\beta^0 = \mu_2 - \mu_1.$$

Note that this is an assumption on the underlying population quantities and will not have practical implications. However, it allows us to have a unique well-defined variable  $\beta^0$  to study the convergence properties of the proposed model. For a discussion on the case  $||L_C^{-1}(\mu_2 - \mu_1)||_{\mathcal{L}^2(\mathcal{M})} = \infty$ , which is related to the perfect classification phenomenon, see Delaigle and Hall (2012), Berrendero et al. (2018), Chen and Jiang (2018), and Kraus and Stefanucci (2019).

The function  $\beta^0$  can be understood as a functional analog of the multivariate discriminant vector of a linear discriminant analysis (Shin 2008). For a new observation with predictor  $x^* \in \mathcal{L}^2(\mathcal{M})$ , it can be used to predict the associated label  $g^*$  with the linear classification rule  $\langle \beta^0, x^* \rangle_{\mathcal{L}^2(\mathcal{M})} > c^{\text{th}}$ , with  $c^{\text{th}}$  an appropriately chosen threshold. Moreover, if X is a Gaussian random function within each group in G, it can be shown that the function  $\beta^0$  defines the linear classifier that minimizes the misclassification error rate, and it is therefore optimal (Delaigle and Hall 2012). The discriminant direction  $\beta^0$  can also be equivalently defined as the minimizer of the functional

$$\frac{1}{2}\langle \beta, L_C \beta \rangle_{\mathcal{L}^2(\mathcal{M})} - \langle \mu_2 - \mu_1, \beta \rangle_{\mathcal{L}^2(\mathcal{M})}. \tag{7}$$

In practice, the population quantities C,  $\mu_1$ , and  $\mu_2$  are unknown and need to be estimated from the data. The goal of a classification model is therefore to recover  $\beta^0$  from the training sample  $\{(g_i, x_i) : i = 1, ..., n\}$  of n independent copies of (G, X). This can be achieved by using the sample covariance  $L_{\hat{C}}$ , with

$$\hat{C}(p,q) = \frac{1}{n} \sum_{i=1}^{n} x_i(p) x_i(q), \qquad p, q \in \mathcal{M}$$

and the sample conditional means  $\hat{\mu}_1$  and  $\hat{\mu}_2$  to replace the population counterparts in equation (7). An estimate  $\hat{\beta}$  of  $\beta^0$  can then be defined as a minimizer of

$$\frac{1}{2}\langle \beta, L_{\hat{C}}\beta \rangle_{\mathcal{L}^2(\mathcal{M})} - \langle \hat{\mu}_2 - \hat{\mu}_1, \beta \rangle_{\mathcal{L}^2(\mathcal{M})} + \mathcal{P}(\beta), \tag{8}$$

where a penalty term  $\mathcal{P}(\beta)$  is typically added to overcome the ill-posedness of the minimization problem, which is due to the low-rank structure of  $L_{\hat{C}}$ . For instance, Park et al. (2021) define a penalty  $\mathcal{P}$  that encourages the estimate  $\hat{\beta}$  to be smooth and sparse, while Kraus and Stefanucci (2019) consider a ridge-type penalty.

The functional discriminant model in equation (8) requires precomputing the empirical covariance function, which is generally not possible for dense functional data supported on multidimensional domains and is ultimately not feasible in our application setting. We therefore propose a direct regularized approach to estimating  $\beta^0$ . This will be possible thanks to the following simple observation. As noted for instance in Delaigle and Hall (2012), the discriminant direction  $\beta^0$  can be equivalently characterized as the solution to the minimization problem

$$\beta^{0} = \operatorname*{arg\,min}_{\beta} \mathbb{E} \left[ Y - \langle X, \beta \rangle_{\mathcal{L}^{2}(\mathcal{M})} \right]^{2}, \tag{9}$$

where Y is an auxiliary scalar random variable such that  $Y = -\frac{1}{\pi_1}$  if G = 1 and  $Y = \frac{1}{\pi_2}$  otherwise. In other words, the classification problem considered can be reformulated as a functional regression problem. This motivates the adoption of a least-squares approach to estimating  $\beta^0$ , based on the empirical counterpart of the objective function in equation (9), where an additional differential regularization term is introduced to incorporate information on the geometric domain  $\mathcal{M}$  and overcome the ill-posedness of the problem. For multivariate data, analogous least-squares formulations have also been adopted, for instance, in Hastie et al. (1994), Mai et al. (2012), and Gaynanova (2020).

### 3.1 Regularized estimation

Given the training sample  $(g_i, x_i)$ , introduce a scalar variable  $y_i$  such that  $y_i = -\frac{n}{n_1}$  if  $g_i = 1$  and  $y_i = \frac{n}{n_2}$  otherwise, where  $n_1$  and  $n_2$  represent the sample sizes of class 1 and 2, respectively. Observe that  $-\frac{n}{n_1}$  and  $\frac{n}{n_2}$  are estimates of the values that can be taken by the random variable Y. Let  $\mathcal{W}^2(\mathcal{M})$  be the Sobolev space of functions in  $\mathcal{L}^2(\mathcal{M})$  with first and second distributional derivatives in  $\mathcal{L}^2(\mathcal{M})$ . We define an estimate  $\hat{\beta} \in \mathcal{W}^2(\mathcal{M})$  of the population quantity  $\beta^0$  as the solution to the following minimization problem

$$\hat{\beta} = \underset{\beta \in \mathcal{W}^2(\mathcal{M})}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^n \left( y_i - \langle x_i, \beta \rangle_{\mathcal{L}^2(\mathcal{M})} \right)^2 + \lambda J(\beta), \tag{10}$$

where the first term is a least-squares estimate of the objective function in equation (9) and the second term is a differential regularization term. The parameter  $\lambda$  controls the trade-off between the least-squares term of the objective function and the penalty term. Our choice of the regularization term is

$$J(\beta) = \|\Delta_{\mathcal{M}}\beta\|_{\mathcal{L}^2(\mathcal{M})}^2 + \varepsilon \|\beta\|_{\mathcal{L}^2(\mathcal{M})}^2, \tag{11}$$

with  $\varepsilon \geq 0$ . This is a linear combination of two terms. The first one is based on the Laplace-Beltrami operator  $\Delta_{\mathcal{M}}: \mathcal{W}^2 \subset \mathcal{L}^2(\mathcal{M}) \to \mathcal{L}^2(\mathcal{M})$  and quantifies the smoothness of the function  $\beta: \mathcal{M} \to \mathbb{R}$  on the nonlinear manifold domain  $\mathcal{M}$ . Specifically, it allows the model estimate  $\hat{\beta}$ , at any fixed point  $p \in \mathcal{M}$ , to borrow strength from the other points on  $\mathcal{M}$  while constraining the 'information' to propagate coherently with the nonlinear manifold structure of the anatomical object  $\mathcal{M}$ . The second term is a generic shrinkage-type regularization.

It is worth noting that the function space  $\mathcal{L}^2(\mathcal{M})$  is linear, even though each function  $f \in \mathcal{L}^2(\mathcal{M})$  is supported on a nonlinear domain. For Euclidean domains, it is common to define a smooth subspace of  $\mathcal{L}^2(\mathcal{M})$  by forming an RKHS from a positive-definite kernel. However, constructing a positive-definite kernel that is compatible with the geometry of a manifold is a non-trivial task (Feragen et al. 2015; Jayasumana et al. 2015). To overcome this challenge, we constructively define a Sobolev norm  $J^{1/2}(\cdot)$  and an associated Sobolev space  $\mathcal{W}^2(\mathcal{M})$  by leveraging a local differential operator, namely the Laplace-Beltrami operator. The discrete counterpart of this local operator is a sparse matrix, reducing our problem to sparse linear algebra and enabling us to solve equation (10) for the large data of our final application. We provide more details about the relationship of our approach with the RKHS approach in Section 4.2.

#### 3.2 Theory

The aim of this section is to provide theoretical guarantees for the performance of the proposed model. Specifically, we provide a probability bound for the out-of-sample risk, i.e., the random variable

$$\mathbb{E}^* \left[ \langle X^*, \beta^0 - \hat{\beta} \rangle_{\mathcal{L}^2(\mathcal{M})} \right]^2, \tag{12}$$

where  $X^*$  is a copy of X that is independent of the training data and  $\mathbb{E}^*$  is the expectation taken over  $X^*$ . Equation (12) measures the discrepancy between the prediction made with the estimated parameter  $\hat{\beta}$  and the 'optimal' prediction made with the unknown population quantity  $\beta^0$ .

Assume for simplicity that  $\varepsilon > 0$ . Then, thanks to the Sobolev embedding theorem (Brezis 2011),  $\exists M \geq 0$  such that for any  $p \in \mathcal{M}$ 

$$f(p) \le \sup_{q} |f(q)| \le M \left( \|\Delta_{\mathcal{M}} f\|_{\mathcal{L}^{2}(\mathcal{M})}^{2} + \varepsilon \|f\|_{\mathcal{L}^{2}(\mathcal{M})}^{2} \right)^{1/2}, \quad \forall f \in \mathcal{W}^{2}(\mathcal{M}),$$

that is, the evaluation operator is a continuous functional. A direct consequence is that the space  $\mathcal{W}^2(\mathcal{M})$  equipped with the norm  $J^{1/2}(\cdot) = \left(\|\Delta_{\mathcal{M}}\cdot\|_{\mathcal{L}^2(\mathcal{M})}^2 + \varepsilon\|\cdot\|_{\mathcal{L}^2(\mathcal{M})}^2\right)^{1/2}$  is an RKHS with a symmetric, positive definite kernel function  $K_{\mathcal{M}}: \mathcal{M} \times \mathcal{M} \to \mathbb{R}$  (Berlinet and Thomas-Agnan 2004). The kernel function  $K_{\mathcal{M}}$  is used only for theoretical investigation here and obtaining its explicit form is, in general, not computationally feasible and not necessary. We will, however, take advantage of the fact that  $L_{K_{\mathcal{M}}}^{1/2}(\mathcal{L}^2(\mathcal{M})) = \mathcal{W}^2(\mathcal{M})$ , where  $L_{K_{\mathcal{M}}}^{1/2}$  denotes the square root of  $L_{K_{\mathcal{M}}}$  (Cucker and Smale 2002). For  $\varepsilon = 0$ , the functional  $J^{1/2}$  defines a semi-norm rather than a norm and similar arguments hold by restricting ourselves to the subspace of  $\mathcal{L}^2(\mathcal{M})$  that is orthogonal to the null space of  $J^{1/2}$ .

Next, we define the sandwich operator  $T=L_{K_M}^{1/2}L_CL_{K_M}^{1/2}$  (Cai and Yuan 2012) and make the following assumptions.

**Assumption 3.1.** The constant  $\kappa^2$ , defined as  $\kappa^2 = \operatorname{ess\,sup} \|L_{K_{\mathcal{M}}}^{1/2} X\|_{\mathcal{L}^2(\mathcal{M})}^2$  is finite.

**Assumption 3.2.** There exists a smooth function  $\beta^0 \in \mathcal{W}^2(\mathcal{M})$  such that  $\beta^0 = L_C^{-1}(\mu_2 - \mu_1)$ .

**Assumption 3.3.** The effective dimension of T satisfies  $D(\lambda) = \text{Tr}((T+\lambda I)^{-1}T) \le c\lambda^{-\theta}$  for constants  $c, \theta > 0$ . Here Tr denotes the trace operator.

Assumption 3.1 allows us to use a Hoeffding-type inequality for Hilbert space valued random elements and has no practical implications. This condition is met, for example, when  $\|X\|_{\mathcal{L}^2(\mathcal{M})}$  is bounded. However, it is more general given that  $L_{K_{\mathcal{M}}}^{1/2}X$  represents a smoothed version of X. Assumption 3.2 guarantees that the population quantity  $\beta^0$  is well-defined and belongs to the space of smooth functions  $\mathcal{W}^2(\mathcal{M})$ . Assumption 3.3 is expressed in terms of properties of the effective dimension  $D(\cdot)$ . For our final choice of  $\lambda$ , it is straightforward to check that this assumption holds by assuming that the eigenvalues  $\{\tau_k\}$  of T decay as  $\tau_k \approx k^{-2r}$ , with  $r > \frac{1}{2}$ . This is a typical assumption in the literature on functional linear models (Cai and Yuan 2012) and is related to the rate of decay of the eigenvalues of  $L_{K_{\mathcal{M}}}$  and  $L_C$ , and their alignment.

The following theorem provides an upper bound for the out-of-sample risk.

**Theorem 3.1.** Under Assumptions 3.1-3.3, if  $\lambda \approx n^{-\frac{1}{1+\theta}}$ , the estimator  $\hat{\beta}$  in equation (10) is such that

$$\mathbb{E}^* \left[ \langle X^*, \beta^0 - \hat{\beta} \rangle_{\mathcal{L}^2(\mathcal{M})} \right]^2 = \mathcal{O}_p \left( n^{-\frac{1}{1+\theta}} \right). \tag{13}$$

Similar rates of convergence have been shown to hold for regularized estimates in the functional linear regression setting (see, e.g., Cai and Yuan 2012; Tong and Ng 2018; Sun et al. 2018; Reimherr et al. 2018). However, a key difference in our model is that the residual random variable  $\varepsilon = Y - \langle X, \beta^0 \rangle_{\mathcal{L}^2(\mathcal{M})}$  and the functional predictor X are not independent, which prevents the direct application of such results. Therefore, Theorem 3.1 shows that in spite of such a dependence structure we are nevertheless able to recover the functional linear model rates of convergence. The proof is provided in Appendix C.

#### 3.3 Nonlinear extensions

To incorporate nonlinearity into the model described in equation (10), one can substitute the term  $\langle x_i, \beta \rangle_{\mathcal{L}^2(\mathcal{M})}$  with a nonlinear function of the data, such as a polynomial term or a single-index model, as done in the context of functional regression models in Yao and Müller (2010) and Jiang and Wang

(2011), respectively. However, these extensions come at the cost of estimating additional functional parameters or optimizing a more complex objective function.

If the covariance structures of the two classes are believed to be different, the proposed functional linear discriminant model can be generalized to an approximate functional quadratic discriminant model, following the approach proposed by Gaynanova and Wang (2019), as follows. We estimate the discriminant rule by minimizing the following objective function with respect to  $\beta_1, \beta_2 \in \mathcal{L}^2(\mathcal{M})$ :

$$\frac{1}{n_1} \sum_{i|g_i=1} \left( 1 - \langle x_i, \beta_1 \rangle_{\mathcal{L}^2(\mathcal{M})} \right)^2 + \frac{1}{n_2} \sum_{i|g_i=2} \left( 1 + \langle x_i, \beta_2 \rangle_{\mathcal{L}^2(\mathcal{M})} \right)^2 + \lambda_1 J(\beta_1) + \lambda_2 J(\beta_2),$$

where  $\lambda_1, \lambda_2 > 0$  are tuning parameters. A modified version of Fisher's criterion (Gaynanova and Wang 2019) is then employed to assign the class by first projecting the data along the estimated directions.

As expected, the simulations presented in Appendix B demonstrate that the approximate functional quadratic discriminant model outperforms the functional linear discriminant model when the covariance structures of the two classes differ. Examining the theoretical properties of this extension is beyond the scope of this paper and is left to future work.

### 4 Additive multivariate generalizations

We now consider a bivariate extension of the functional model introduced in Section 3, which incorporates the geometric component of the original data. We therefore consider the training sample  $\{(g_i, v_i, x_i)\}$ , where  $v_i \in \mathcal{V}(\mathbb{R}^3)$  is a vector field representing the subject-specific geometry of the *i*th subject. Recall that  $(\mathcal{V}(\mathbb{R}^3), \|\cdot\|_{\mathcal{V}(\mathbb{R}^3)})$  is an RKHS of smooth functions with compact support in  $\mathbb{R}^3$ . Moreover, denote by  $K_{\mathbb{R}^3}$  its reproducing kernel.

Suppose the training set  $\{(g_i, v_i, x_i)\}$  consists of n independent copies of (G, V, X), a triplet of random variables with V a zero-mean random function taking values in  $\mathcal{V}(\mathbb{R}^3)$ , X a zero-mean random function taking values in  $\mathcal{L}^2(\mathcal{M})$ , and G a binary random variable such that  $P(G = 1) = \pi_1$  and  $P(G = 2) = \pi_2$ .

We now adopt the multivariate functional data notation from Happ and Greven (2018), and define  $X(p) = (V(p_1), X(p_2))$ , with  $p = (p_1, p_2) \in D = D_1 \times D_2 = \mathbb{R}^3 \times \mathcal{M}$ . The multivariate random function X takes values in a Hilbert space  $\mathcal{H} = \mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3) \times \mathcal{L}^2(\mathcal{M})$  with inner product  $\langle f, g \rangle_{\mathcal{H}} = \langle f^{(1)}, g^{(1)} \rangle_{\mathcal{L}^2(\mathbb{R}^3, \mathbb{R}^3)} + \langle f^{(2)}, g^{(2)} \rangle_{\mathcal{L}^2(\mathcal{M})}$  for  $f, g \in \mathcal{H}$ . Here  $f^{(j)}$ , with  $j \in \{1, 2\}$ , denotes the jth functional component of the multivariate function f. For  $p, q \in D$ , define the matrix of covariances  $C(p, q) = \mathbb{E}(X(p) \otimes X(q))$  with elements  $C_{lj}(p_l, q_j) = \mathbb{E}[X^{(l)}(p_l)X^{(j)}(q_j)]$  where  $p_l \in D_l, q_j \in D_j, l \in \{1, 2\}$ , and  $j \in \{1, 2\}$ . Denote the conditional means of X by  $\mu_1 = \left(\mu_1^{(1)}, \mu_1^{(2)}\right) := (\mathbb{E}[V|G=1], \mathbb{E}[X|G=1])$  and  $\mu_2 = \left(\mu_2^{(1)}, \mu_2^{(2)}\right) := (\mathbb{E}[V|G=2], \mathbb{E}[X|G=2])$ , and assume  $\mu_1 \neq \mu_2$ . The covariance operator  $L_C : \mathcal{H} \to \mathcal{H}$  is such that the jth component of  $L_C f$ , for any  $f \in \mathcal{H}$ , is given by

$$(L_{\mathbf{C}}\mathbf{f})^{(j)}(p_j) = \sum_{i=1}^{2} \int_{D_i} C_{ij}(q_i, p_j) f^{(i)}(q_i) dq_i.$$
(14)

Similar to the univariate case, we assume that the population quantity  $\beta^0 \in \mathcal{H}$  is well-defined and satisfies the equation

$$L_{\mathbf{C}}\boldsymbol{\beta}^0 = \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1.$$

This can be viewed as a multivariate generalization of the linear discriminant direction defined in the previous section. We now turn to the problem of defining an estimator for  $\beta^0$ .

Estimate	Function space	Norm	Kernel		
$\hat{eta}^F:\mathcal{M} o\mathbb{R}$	$\mathcal{W}^2(\mathcal{M})$	$J^{1/2}(\cdot) = \left(\ \Delta_{\mathcal{M}} \cdot\ _{\mathcal{L}^2(\mathcal{M})}^2 + \varepsilon\ \cdot\ _{\mathcal{L}^2(\mathcal{M})}^2\right)^{1/2}$	$K_{\mathcal{M}}$ (implicit)		
$\hat{\beta}^G: \mathbb{R}^3 \to \mathbb{R}^3$	$\mathcal{V}(\mathbb{R}^3)$	$\ \cdot\ _{\mathcal{V}(\mathbb{R}^3)}  ext{ (implicit)}$	$K_{\mathbb{R}^3}$		

Table 1: Table summarizing estimates and associated function spaces, norms and kernels.

#### 4.1 Regularized estimation

Let the variable  $y_i$  be such that  $y_i = -\frac{n}{n_1}$  if  $g_i = 1$  and  $y_i = \frac{n}{n_2}$  otherwise. We define the multivariate functional estimate  $\hat{\beta} = (\hat{\beta}^G, \hat{\beta}^F)$  of the population quantity  $\beta^0$  to be the solution to the following minimization problem

$$\left(\hat{\beta}^{G}, \hat{\beta}^{F}\right) = \underset{\substack{\beta^{G} \in \mathcal{V}(\mathbb{R}^{3}) \\ \beta^{F} \in \mathcal{W}^{2}(\mathcal{M})}}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^{n} \left(y_{i} - \langle v_{i}, \beta^{G} \rangle_{\mathcal{L}^{2}(\mathbb{R}^{3}, \mathbb{R}^{3})} - \langle x_{i}, \beta^{F} \rangle_{\mathcal{L}^{2}(\mathcal{M})}\right)^{2} + \lambda_{1} \|\beta^{G}\|_{\mathcal{V}(\mathbb{R}^{3})}^{2} + \lambda_{2} J(\beta^{F}), (15)$$

with  $\lambda_1, \lambda_2$  tuning parameters.

Equation (15) extends the model proposed in Section 3 to account for both the geometric functional descriptor  $v_i$  and the function  $x_i$  in an additive fashion. The regularization terms in the equation enforce smoothness on the functional estimates  $\hat{\beta}^G$  and  $\hat{\beta}^F$  in their respective function spaces.

#### 4.2 Differential regularization and kernel penalty: a unified modeling perspective

In Section 4.1, we have adopted two different approaches to produce smooth estimates  $\hat{\beta}^G : \mathbb{R}^3 \to \mathbb{R}^3$  and  $\hat{\beta}^F : \mathcal{M} \to \mathbb{R}$ . The smoothness of  $\hat{\beta}^F$  is enforced by means of a penalty  $J(\cdot)$  defined in terms of a Sobolev norm, which implicitly defines a kernel  $K_{\mathcal{M}}$ . On the other hand, the smoothness of  $\hat{\beta}^G$  is enforced by means of a norm  $\|\cdot\|_{\mathcal{V}(\mathbb{R}^3)}$ , defined implicitly through the direct definition of a kernel  $K_{\mathbb{R}^3}$ . For clarity, we summarize the relevant function spaces and associated norms and kernels in Table 1.

From a methodological perspective, the reproducing kernel  $K_{\mathbb{R}^3}(p,q)$  can be understood as a measure of the influence of the function value at  $p \in \mathbb{R}^3$  on the function value at  $q \in \mathbb{R}^3$  and vice-versa. Defining a smooth function space through a kernel has arguably an advantage when it comes to discretizing an infinite-dimensional minimization problem over that function space. In fact, thanks to the well-known representer theorem (Wahba 1990; Yuan and Cai 2010), under mild conditions, its exact solution can be expressed as a linear combination of the elements of a n-dimensional basis, which involves the explicit expression of the kernel.

Hence, it is natural to wonder whether a similar approach could be adopted for  $\hat{\beta}^F: \mathcal{M} \to \mathbb{R}$ . In other words, can we define a smooth real function space on  $\mathcal{M}$  by explicitly defining a kernel  $K_{\mathcal{M}}: \mathcal{M} \times \mathcal{M} \to \mathbb{R}$  encoding a measure of influence that is coherent with the nonlinear geometry of  $\mathcal{M}$ ? This, however, is a challenging task due to the positive-definiteness property that  $K_{\mathcal{M}}$  must satisfy. Consider, for instance, the popular exponential kernel. Its natural extension to the manifold setting is  $K_{\mathcal{M}}(p,q) = \exp(-c \, \mathrm{d}_{\mathcal{M}}(p,q)^2)$ , where  $\mathrm{d}_{\mathcal{M}}(p,q)$  is the geodesic distance between  $p \in \mathcal{M}$  and  $q \in \mathcal{M}$ . Unfortunately, this kernel cannot be guaranteed to be positive definite for a general nonlinear manifold  $\mathcal{M}$  (Feragen et al. 2015; Jayasumana et al. 2015).

Alternatively, we could try to compute an explicit form of the kernel  $K_{\mathcal{M}}$  from  $J^{1/2}(\cdot)$  by employing the following identity (Wahba 1990; Fasshauer and Ye 2013)

$$f(p) = \langle K_{\mathcal{M}}(p,\cdot), f \rangle_{\mathcal{W}^2(\mathcal{M})}, \qquad \forall p \in \mathcal{M}, f \in \mathcal{W}^2(\mathcal{M}),$$
 (16)

where  $\langle \cdot, \cdot \rangle_{\mathcal{W}^2(\mathcal{M})}$  is the inner product that induces the norm  $J^{1/2}(\cdot)$ . However, closed-form solutions to equation (16) are not available in our setting. Approximate solutions could be computed by Finite Elements Analysis (FEA) (Quarteroni 2009), but we would still face the challenge of storing the dense

object  $K_{\mathcal{M}}(\cdot,\cdot)$ . As described in Section 5, we instead leverage FEA to directly discretize the function  $\beta^F: \mathcal{M} \to \mathbb{R}$  in equation (15).

This highlights that the choice of the two modeling approaches is not arbitrary and that, arguably, for functional estimates supported on Euclidean spaces, defining explicitly a reproducing kernel is likely the preferred choice. Meanwhile, for general non-Euclidean domains, where defining a reproducing kernel is not trivial, the differential penalization approach is preferable.

#### 4.3 Theory

Define the diagonal matrix of reproducing kernels K(p,q) with entries  $K_{11}(p_1,q_1)=K_{\mathbb{R}^3}(p_1,q_1)$  and  $K_{22}(p_2,q_2)=K_{\mathcal{M}}(p_2,q_2); p_i\in D_i, p_j\in D_j.$  Let  $L_{\mathbf{K}}:\mathcal{H}\to\mathcal{H}$  be the associated integral operator and, analogously to the univariate functional setting, define the sandwich operator  $T=L_{\pmb{K}}^{1/2}L_{\pmb{C}}L_{\pmb{K}}^{1/2}$ . We make the following assumptions, which are analogous to Assumptions 3.1-3.3.

**Assumption 4.1.** The constant  $\kappa_2^2$ , defined as  $\kappa_2^2 = \text{ess sup } \|L_K^{1/2} X\|_{\mathcal{H}}^2$  is finite.

**Assumption 4.2.** There exists a smooth function  $\beta^0 \in \mathcal{V}(\mathbb{R}^3) \times \mathcal{W}^2(\mathcal{M})$  such that  $\beta^0 = L_{\mathbf{C}}^{-1}(\mu_2 - \mu_1)$ .

**Assumption 4.3.** The penalty coefficient  $\lambda := \lambda_1 = \lambda_2$  and the effective dimension of T satisfy  $D(\lambda) = \text{Tr}((T + \lambda I)^{-1}T) \le c\lambda^{-\theta} \text{ for constants } c, \theta > 0.$ 

The following theorem, which is an extension of Theorem 3.1, provides an upper bound for the out-of-sample risk.

**Theorem 4.1.** Under Assumptions 4.1-4.3, if  $\lambda \approx n^{-\frac{1}{1+\theta}}$ , the estimator  $\hat{\beta} = (\hat{\beta}^G, \hat{\beta}^F)$  in equation (15) is such that

$$\mathbb{E}^* \left[ \langle \boldsymbol{X}^*, \boldsymbol{\beta}^0 - \hat{\boldsymbol{\beta}} \rangle_{\mathcal{H}} \right]^2 = \mathcal{O}_p \left( n^{-\frac{1}{1+\theta}} \right), \tag{17}$$

where  $X^*$  is a copy of X that is independent of the training data and  $\mathbb{E}^*$  is the expectation taken over  $X^*$ .

#### 5 Discretization

Consider a triangle mesh, denoted by  $\mathcal{M}_{\mathcal{T}}$ , which is formed by the union of a finite set of triangles,  $\mathcal{T}$ . These triangles share a common set of s vertices, denoted as  $\xi_1, \ldots, \xi_s$ . Let  $\mathcal{M}_{\mathcal{T}}$  be an approximate representation of the manifold  $\mathcal{M}$ . We then introduce the linear finite element space  $\mathcal{W}_{\mathcal{T}}$  consisting of a set of globally continuous functions over  $\mathcal{M}_{\mathcal{T}}$  that are affine within each triangle  $\tau$  in  $\mathcal{T}$ , i.e.,

$$\mathcal{W}_{\mathcal{T}} = \{ w \in C^0(\mathcal{M}_{\mathcal{T}}) : w|_{\tau} \text{ is affine } \forall \tau \in \mathcal{T} \}.$$

The space  $\mathcal{W}_{\mathcal{T}}$  is spanned by the Finite Elements (FE) basis  $\psi_1, \ldots, \psi_s$ , where  $\psi_l(\xi_i) = 1$ , if l = j, and  $\psi_l(\xi_j) = 0$  otherwise. In Figure 2, we show one element of this basis. Moreover, define  $\psi$  as the vector-valued function  $\psi(\cdot) = (\psi_1(\cdot), \dots, \psi_s(\cdot))^T$ . Our goal is to find an approximate solution  $\hat{\beta}_{\mathcal{T}}^F$  of the form

$$\beta_{\mathcal{T}}^{F}(\cdot) = \sum_{l=1}^{s} c_{l}^{F} \psi_{l}(\cdot) = (\boldsymbol{c}^{F})^{T} \boldsymbol{\psi}(\cdot), \tag{18}$$

where  $\mathbf{c}^F = \left(c_1^F, \cdots, c_s^F\right)^T$ . Let now M and S be the sparse mass and stiffness  $s \times s$  matrices defined as  $(M)_{jj'} = \int_{\mathcal{M}_T} \psi_j \psi_{j'}$ and  $(S)_{jj'} = \int_{\mathcal{M}_{\mathcal{T}}} \nabla_{\mathcal{M}_{\mathcal{T}}} \psi_j \cdot \nabla_{\mathcal{M}_{\mathcal{T}}} \psi_{j'}$ , where  $\nabla_{\mathcal{M}_{\mathcal{T}}}$  is the gradient operator on the mesh  $\mathcal{M}_{\mathcal{T}}$ . For  $\beta_{\mathcal{T}}^F$  of the form given in equation (18), we have that the penalty term  $J(\cdot)$  can be approximated as  $(c^F)^T D_{\mathcal{M}_{\mathcal{T}}} c^F$ , with  $D_{\mathcal{M}_{\mathcal{T}}} = SM^{-1}S + \varepsilon M$  (Lila et al. 2016). Further, following an approach often adopted in FEA (Fried and Malkus 1975; Hinton et al. 1976), we replace the dense matrix  $M^{-1}$  with the sparse matrix  $\tilde{M}^{-1}$ , where  $\tilde{M}$  is the diagonal matrix such that  $\tilde{M}_{jj} = \sum_{j'} M_{jj'}$ . This results in the sparse penalty matrix  $D_{\mathcal{M}_{\mathcal{T}}} = S\tilde{M}^{-1}S + \varepsilon M$ . In practice, each functional observation  $x_i$  is also of the form given in equation (18). Therefore, denoting by  $\mathbb{X}$  the  $n \times s$  matrix where each row consists of the basis coefficients of  $x_i$ , the terms  $\{\langle x_i, \beta^F \rangle_{\mathcal{L}^2(\mathcal{M})}\}$  can be approximated by the entries of the vector  $\mathbb{X}M\mathbf{c}^F$ .

We now turn our attention to the estimate  $\hat{\beta}^G \in \mathcal{V}(\mathbb{R}^3)$ . Since for  $\mathcal{V}(\mathbb{R}^3)$  we have an explicit form of the associated reproducing kernel  $K_{\mathbb{R}^3}$ , we employ the representer theorem (Wahba 1990; Yuan and Cai 2010) and take  $\hat{\beta}^G$  of the form

$$\beta^G(\cdot) = \sum_{i=1}^n c_i^G \int_{\mathbb{R}^3} K_{\mathbb{R}^3}(p, \cdot) v_i(p) \, dp. \tag{19}$$

On the right hand side of Figure 2, we show an example of a basis function  $\int_{\mathbb{R}^3} K_{\mathbb{R}^3}(p,\cdot)v_i(p) dp$ . Then, we have that  $\|\beta^G\|_{\mathcal{V}}^2 = (\mathbf{c}^G)^T \Sigma \mathbf{c}^G$ , where  $\mathbf{c}^G = \left(c_1^G, \dots, c_n^G\right)^T$  and  $\Sigma$  is a  $n \times n$  matrix with entries

$$\Sigma_{ij} = \int \int v_i(p)^T K_{\mathbb{R}^3}(p,q) v_j(q) \, dp \, dq.$$

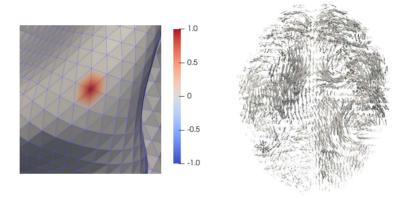


Figure 2: On the left hand side, we show an element of the FE basis  $\{\psi_l : \mathcal{M}_{\mathcal{T}} \to \mathbb{R}, l = 1, \dots, s\}$ . This is a scalar affine function within each triangle of the mesh  $\mathcal{M}_{\mathcal{T}}$  that takes value 1 on a fixed vertex and value 0 on every other vertex. On the right hand side, we show an element of the basis  $\{\int_{\mathbb{R}^3} K_{\mathbb{R}^3}(p,\cdot)v_i(p)\,dp, i=1,\dots,n\}$ . This is a smooth vector-valued function from  $\mathbb{R}^3$  to  $\mathbb{R}^3$ .

As a result, the coefficients of the approximate solution of the model in equation (15) are given by

$$(\hat{\boldsymbol{c}}^G, \hat{\boldsymbol{c}}^F) = \underset{\boldsymbol{c}^G \in \mathbb{R}^n, \boldsymbol{c}^F \in \mathbb{R}^s}{\min} \|\boldsymbol{y} - \boldsymbol{\Sigma} \boldsymbol{c}^G - \boldsymbol{X} \boldsymbol{M} \boldsymbol{c}^F\|_2^2 + \lambda_1 (\boldsymbol{c}^G)^T \boldsymbol{\Sigma} \boldsymbol{c}^G + \lambda_2 (\boldsymbol{c}^F)^T D_{\mathcal{M}_{\mathcal{T}}} \boldsymbol{c}^F,$$
(20)

where  $\mathbf{y} = (y_1, \dots, y_n)^T$  is the vector of auxiliary response variables. It is easy to check that this minimization problem can be equivalently written as the following augmented quadratic least-squares problem

$$\left(\hat{\boldsymbol{c}}^{G}, \hat{\boldsymbol{c}}^{F}\right) = \underset{\boldsymbol{c}^{G} \in \mathbb{R}^{n}, \boldsymbol{c}^{F} \in \mathbb{R}^{s}}{\arg\min} \left\| \begin{bmatrix} \boldsymbol{y} \\ \boldsymbol{0} \end{bmatrix} - A \begin{bmatrix} \boldsymbol{c}^{G} \\ \boldsymbol{c}^{F} \end{bmatrix} \right\|_{2}^{2}, \tag{21}$$

where **0** is the zero-vector of length 2s + n and with

$$A = egin{bmatrix} \Sigma & \mathbb{X}M \ 0 & \lambda_2^{rac{1}{2}} ilde{M}^{-rac{1}{2}}S \ 0 & \lambda_2^{rac{1}{2}}arepsilon^{rac{1}{2}}M^{rac{1}{2}} \ \lambda_1^{rac{1}{2}}\Sigma^{rac{1}{2}} & 0 \end{bmatrix}.$$

Note that for  $n \ll s$ , which is the setting of our application, the matrix A is sparse. Therefore, the minimization problem (21) can be efficiently solved by conjugate gradients, or its variations, e.g. LSQR (Paige and Saunders 1982), without requiring the explicit computation of the high-dimensional normal matrix  $A^TA$  – a quantity related to the covariance structure of the functional predictors.

An approximate solution to the univariate model in equation (10) follows as a special case of the multivariate case considered here.

### 6 Application

#### 6.1 Data and preprocessing

We analyze a cohort of n=484 subjects from the ADNI and PPMI studies. On the basis of the ATN classification scheme (Jack et al. 2016), each subject was assigned to one of the two diagnostic categories – C: Control ( $n_1=100$ ) and AD: Alzheimer's Disease ( $n_2=384$ ). Here, we focus on data collected at the baseline visit, which includes, among other imaging modalities, structural T1-weighted MRI.

The T1-weighted images were preprocessed using FreeSurfer (Dale et al. 1999; Fischl et al. 1999). Specifically, white matter, grey matter, and cerebrospinal fluid were segmented and used to extract the outer and inner surfaces of the cerebral cortex. From these two surfaces, we generated a central surface interpolating the midpoints between the outer and inner layers, which offers an accurate representation of the two-dimensional anatomical structure of the cerebral cortex. This representation has the benefit of encoding a notion of distance between brain regions that is neurologically more relevant than the original volumetric representation. The cortical surface can moreover be coupled with one or more maps describing complementary structural or functional properties of the cortex, such as cortical thickness measurements (Fischl and Dale 2000), fMRI signals, or connectivity maps (Smith et al. 2013; Yeo et al. 2011). In this study, we focus on cortical thickness, which is estimated from the distances between the outer and inner surfaces of the cerebral cortex. Next, the n surfaces were registered and sub-sampled.

As a result of the preprocessing stage, we obtain n=484 triangle meshes  $\{\mathcal{M}_i^{\mathcal{T}}\}$  consisting of s=64K vertices in correspondence across subjects, along with a set of triangular faces defining how these vertices are connected to delineate the surfaces. By classical Generalized Procrustes Analysis (Dryden and Mardia 2016), translation, rigid rotation, and scale were removed from the data, while jointly estimating a template  $\mathcal{M}^{\mathcal{T}}$ , which is also a triangle mesh with s=64K vertices in correspondence with those of the individual subjects. Each surface has been moreover coupled with cortical thickness measurements at the mesh vertices, which we model as a real piecewise linear continuous function  $z_i^{\mathcal{T}}$  over the mesh  $\mathcal{M}_i^{\mathcal{T}}$ .

The preprocessing stage results in a set of FoSs  $\{(\mathcal{M}_i^T, z_i^T), i = 1, ..., n\}$ , which are discretized versions of the continuous idealized objects  $\{(\mathcal{M}_i, z_i), i = 1, ..., n\}$  introduced in Section 2. To simplify the notation, we drop the superscript  $\mathcal{T}$ . Moreover, we denote the diagnostic labels by  $\{g_i \in \{C, AD\}, i = 1, ..., n\}$ . Three examples of such FoSs, and associated diagnostic labels, are given in Figure 1A.

Here, we are interested in using the proposed models to identify subjects with AD from the extracted FoSs. The interpretability of these methods is an important feature. Indeed, while it is crucial to build models with good classification accuracy, it is equally important to describe the estimated relationship between the predictors and the outcome variable, in order to inform subsequent studies and generate data-driven hypotheses about the pathophysiology of AD.

#### 6.2 Linear functional representations

For each FoS, we compute a function  $v_i \in \mathcal{V}(\mathbb{R}^3)$  such that  $\varphi_{v_i}(\mathcal{M})$  closely approximates  $\mathcal{M}_i$ , where closeness is measured as the sum of Euclidean distances between the corresponding vertices of  $\varphi_{v_i}(\mathcal{M})$  and  $\mathcal{M}_i$ . As noted in Appendix A.1, alternative definitions of surface similarity could also be used.

We can then transport the function  $z_i : \mathcal{M}_i \to \mathbb{R}$  onto the template defining a continuous piecewise linear function  $x_i = z_i \circ \varphi_{v_i}$ . This leads to the definition of the bivariate functional representation  $(v_i, x_i)$  that is a linear representation of the FoS  $(\mathcal{M}_i, z_i) \approx (\varphi_{v_i}(\mathcal{M}), x_i \circ \varphi_{v_i}^{-1})$ . Further details on the computation of  $v_i$  can be found in Appendix A.1.

#### 6.3 Discriminant analysis

Our aim is to estimate directions in the parametrizing geometric and thickness spaces that are most predictive of AD. To this end, we apply the model introduced in Section 4.1 to the training data  $\{g_i, v_i - \bar{v}, x_i - \bar{x}\}$ , where  $\bar{v}$  and  $\bar{x}$  are the sample means of  $\{v_i\}$  and  $\{x_i\}$ . The training data are a subset of the entire dataset containing 50% of the observations. From this process, we derive the estimates  $\hat{\beta}^G : \mathbb{R}^3 \to \mathbb{R}^3$  and  $\hat{\beta}^F : \mathcal{M} \to \mathbb{R}$ . Given a new subject with predictors  $(v^*, x^*)$ , these estimates can be used in conjunction with the classification rule  $\langle v^* - \bar{v}, \hat{\beta}^G \rangle + \langle x^* - \bar{x}, \hat{\beta}^F \rangle > c^{\text{th}}$  to predict the diagnostic label of a new subject. The cut-off level  $c^{\text{th}}$  can be chosen by computing sensitivity and specificity on a test set, for different values  $c^{\text{th}}$ , and by selecting the desired level and type of accuracy.

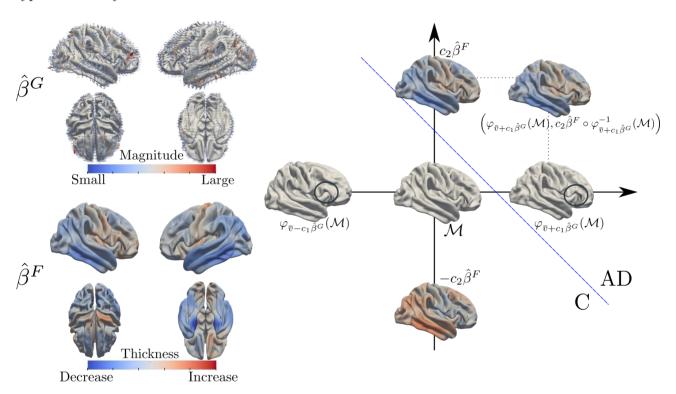


Figure 3: On the left hand side, we show the most discriminant geometric and thickness directions as estimated from the linear representations  $\{(v_i - \bar{v}, x_i - \bar{x})\}$ . These are a vector field  $\hat{\beta}^G : \mathbb{R}^3 \to \mathbb{R}^3$ , representing the most predictive geometric pattern of AD, and a function  $\hat{\beta}^F : \mathcal{M} \to \mathbb{R}$ , representing the most predictive cortical thickness pattern of AD. For a new FoS, with linear representation  $(v^*, x^*)$ , we compute the score  $\langle v^* - \bar{v}, \hat{\beta}^G \rangle + \langle x^* - \bar{x}, \hat{\beta}^F \rangle$  and predict whether the subject has AD by comparing the score value with a predetermined threshold  $c^{\text{th}}$ . On the right hand side, we depict the process of mapping back the estimates  $\hat{\beta}^G$  and  $\hat{\beta}^F$  to the space of FoSs. On the same space, we also pictorially map the classification rule adopted. In the  $\hat{\beta}^F$  figure, the blue regions represent the areas of the cortical surface where a thinner cortex, relative to the population average, is indicative of AD. These are mostly localized in the lateral temporal, entorhinal, inferior parietal, precuneus, and posterior cingulate cortices. The red arrows in the  $\hat{\beta}^G$  figure represent the regions where differences in the morphological configuration of the cerebral cortex, compared to the population average, are most predictive of AD. The specific types of morphological changes can be inspected by comparing the surfaces  $\varphi_{\bar{\nu}-c_1\hat{\beta}G}(\mathcal{M})$  and  $\varphi_{\bar{\nu}+c_1\hat{\beta}G}(\mathcal{M})$ , on the right hand side diagram.

These estimates effectively identify linear directions  $\{c_1\hat{\beta}^G, c_1 \in \mathbb{R}\}$  and  $\{c_2\hat{\beta}^F, c_2 \in \mathbb{R}\}$ , in their respective spaces, that can be interpreted as the most discriminant geometric and thickness directions. Specifically, large values of  $c_1 \in \mathbb{R}$  and  $c_2 \in \mathbb{R}$  describe configurations  $c_1\hat{\beta}^G$  and  $c_2\hat{\beta}^F$  that are predictive of AD. Low values of  $c_1 \in \mathbb{R}$  and  $c_2 \in \mathbb{R}$  describe configurations that are instead predictive of the subject being healthy. Moreover, given the additive modeling assumption on the geometric and thickness components, for every configuration  $(c_1\hat{\beta}^G, c_2\hat{\beta}^F)$ , an increase of  $c_1$  ( $c_2$ ) for a fixed  $c_2$  ( $c_1$ ) describes a configuration that is more strongly associated with AD.

Crucially, these linear trajectories on the parametrizing space can be mapped back to the original space of FoSs by using equation (2), defining the curved space

$$\left(\varphi_{\bar{v}+c_1\hat{\beta}^G}(\mathcal{M}), c_2\hat{\beta}^F \circ \varphi_{\bar{v}+c_1\hat{\beta}^G}^{-1}\right), \qquad c_1, c_2 \in \mathbb{R}.$$

We fit the proposed model for different choices of the parameters  $\lambda_1$  and  $\lambda_2$ . Recall that  $\lambda_1$  controls the regularity of the geometric discriminant direction and  $\lambda_2$  that of the thickness discriminant direction, with high values virtually constraining the solution to be the zero function. The final choice of  $\lambda_1$  and  $\lambda_2$  is the result of a compromise between classification accuracy on the test set and the consistency of the estimated discriminant directions with the neurodegenerative nature of the disease (see also Discussion). The test Area under the ROC Curve (AUC) of the selected model is 0.7006.

#### 6.4 Results

On the left hand side of Figure 3, we show the estimated most discriminant geometric and thickness directions, i.e.,  $\hat{\beta}^G: \mathbb{R}^3 \to \mathbb{R}^3$  and  $\hat{\beta}^F: \mathcal{M} \to \mathbb{R}$ . These have been estimated by applying the model in equation (15) to the linear representations  $\{(v_i - \bar{v}, x_i - \bar{x})\}$ . The colormap describing  $\hat{\beta}^F$  illustrates what types of variations, with respect to the population average cortical thickness, are most predictive of AD. Specifically, a thinner cerebral cortex in the blue areas (i.e., lateral temporal, entorhinal, inferior parietal, precuneus, and posterior cingulate cortices) is associated with AD. These results are consistent with the typical thickness signature of AD observed to date (see, e.g., Bondareff et al. 1989; Dickerson et al. 2009; Sabuncu et al. 2011). The geometric component  $\hat{\beta}^G$  is instead a vector field in  $\mathbb{R}^3$ . This is a linear representation of the morphological variations, with respect to the population average cortex geometry, that are associated with AD. While a full understanding of its meaning is only possible by mapping  $\hat{\beta}^G$  back to the space of FoSs, i.e., by examining  $\varphi_{\bar{v}+c_1\hat{\beta}^G}(\mathcal{M})$  for different choices of  $c_1 \in \mathbb{R}$ , the magnitude of the vector field  $\hat{\beta}^G$ , at any fixed point, offers a rough indication of the cortical regions whose morphological variations are most relevant to the classification problem.

On the right hand side of Figure 3, we show the FoSs associated with the linear representations  $\hat{\beta}^G$  and  $\hat{\beta}^F$ , that is,  $\left(\varphi_{\bar{v}+c_1\hat{\beta}^G}(\mathcal{M}), c_2\hat{\beta}^F \circ \varphi_{\bar{v}+c_1\hat{\beta}^G}^{-1}\right)$  with  $c_1, c_2 \in \mathbb{R}$ . These describe the most predictive patterns of AD in terms of the original neurobiological objects. We have circled a specific area of the brain to ease comparison and highlight the morphological patterns that the model deems relevant to the classification problem.

#### 6.5 Comparison against alternative approaches

In this section, we compare the test AUC of our proposed classification method with alternative models and evaluate different representation models for FoSs. In addition to the functional linear discriminant model (FLDA) that we propose, we also consider the following alternatives: (i) FPCA+LDA: The geometry-aware FPCA model proposed in Lila et al. (2016), followed by multivariate LDA (Hastie et al. 2009) on the PC scores; (ii) Lasso: A logistic regression model with lasso regularization (Tibshirani 1996); (iii) Ridge: A logistic regression model with an  $\ell^2$  regularization (Hoerl and Kennard 1970); (iv) FQDA: The approximate functional quadratic discriminant model defined in Section 3.3; (v) RF: A Random forest model (Breiman 2001); (vi) SVM: A support vector machine with a squared exponential kernel (Cortes and Vapnik 1995); (vii) NN: A multilayer feedforward neural network (Hastie et al. 2009).

Representation model	Linear methods				Nonlinear methods			
	FLDA	FPCA+LDA	Lasso	Ridge	FQDA	RF	SVM	NN
Thickness	0.7626	0.7583	0.7487	0.7632	0.7710	0.6043	0.7597	0.7678
Thickness & Displacement	0.6623	-	0.6626	0.6571	-	0.6742	0.6861	0.6771
Thickness & Shape spectrum	-	-	0.7832	0.6638	-	0.5797	0.7606	0.6878
Proposed FoSs representation	0.7716	-	0.7484	0.7646	-	0.7132	0.7600	0.7443

Table 2: The test AUC of the classification methods applied to the data of our final application. Four different representation models have been considered: (i) the registered thickness map  $x_i : \mathcal{M} \to \mathbb{R}$  without geometric information; (ii) the parametrization  $h_i : \mathcal{M} \to \mathbb{R}^4$ , where the first three components are the surface coordinates, and the last component is thickness; (iii) the registered thickness map  $x_i : \mathcal{M} \to \mathbb{R}$  and the first 200 eigenvalues of the Laplace-Beltrami operator computed on the surface  $\mathcal{M}_i$ ; and (iv) the proposed representation model  $(x_i, v_i)$ . The symbol '-' indicates that although the method could be adapted to accommodate the specific FoS representation model, its implementation is beyond the scope of this paper and is left to future work. For each representation model, the top-performing method is highlighted.

Furthermore, besides the proposed representation model  $(v_i, x_i)$  for FoSs, we also consider the following representations: (i) Thickness: Spatially normalized thickness maps  $x_i : \mathcal{M} \to \mathbb{R}$  without geometric information; (ii) Thickness & Displacement: The parametrizations  $h_i : \mathcal{M} \to \mathbb{R}^4$ , where the first three components are the surface coordinates, and the last component is the (spatially normalized) thickness map; (iii) Thickness & Shape spectrum: The spatially normalized thickness maps  $x_i : \mathcal{M} \to \mathbb{R}$  and the first 200 eigenvalues of the Laplace-Beltrami operator computed on the surface  $\mathcal{M}_i$ , i.e., a spectral representation of shape (Reuter et al. 2006).

To evaluate the listed methods and representation models, we split the dataset into three sets, namely the training set, validation set, and test set, comprising 50%, 20%, and 30% of the data, respectively. While a Monte Carlo evaluation of these methods would be desirable, it is computationally prohibitive, so we defer that analysis to the simulation setting in Appendix B. However, we use the same exact data split for all methods. The models are trained on the training set, hyperparameters are chosen to maximize the AUC on the validation set and the selected model is tested on the test set, resulting in the AUC scores presented in Table 2. Note that, in contrast to the results presented in Section 6.4 and Figure 3, all hyperparameters of the proposed methods have been chosen to maximize the AUC on the validation set, rather than striking a balance between the classification accuracy and consistency of the estimated discriminant directions with the neurodegenerative nature of the disease. Hence, the test AUC value of the proposed method is different from that in Section 6.4.

For standard multivariate models, we use the values of  $x_i$  at the vertices of the template mesh (64K values) and the RKHS coefficients of the estimated  $v_i$  (192K coefficients) to construct the data matrix. We have also implemented a variation of the functional linear discriminant model introduced in Section 4, for multivariate functions whose components share a non-linear domain  $\mathcal{M}$ , in order to accommodate the representation  $h_i: \mathcal{M} \to \mathbb{R}^4$ .

The results are shown in Table 2, from which we can make several observations. Firstly, if the goal is to maximize prediction accuracy, then the best-performing model is a lasso-penalized generalized linear model applied to thickness maps and the first 200 eigenvalues of the Laplace-Beltrami operator of the surfaces. However, as mentioned in the introduction, this "lossy" shape representation cannot be mapped back to the original space of neurobiological objects, leading to a less interpretable model. Additionally, our results show that in the context of our application, the representation  $(v_i, x_i)$  performs better across all methods than using the representation  $h_i : \mathcal{M} \to \mathbb{R}^4$ . The latter appears to be more susceptible to overfitting, resulting in inferior performance even when compared to models that use thickness only. Finally, classification using thickness alone produces satisfactory results. The top-performing models are the proposed FLDA and FQDA, and the ridge logistic regression model. One possible explanation for this is that the registered thickness maps may include some geometric

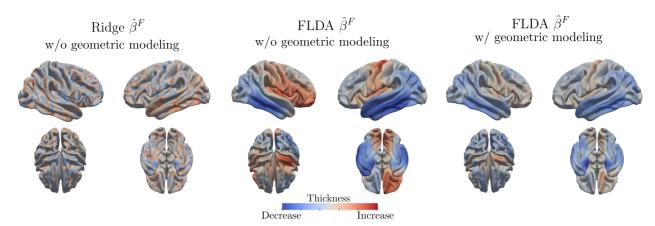


Figure 4: On the left side, we show the discriminant direction derived from applying a ridge logistic regression model to the thickness maps. In the center, we show the discriminant direction resulting from fitting the proposed model in equation (10) to the thickness maps. Although it does not account for subject-specific geometric variations, this model enforces smoothness. On the right side, we have the cortical thickness discriminant direction obtained by fitting the model in equation (15), which explicitly accounts for inter-subject geometric differences. The results of the logistic regression are more difficult to interpret due to the high spatial variability. The model in equation (10) provides more interpretable results thanks to its smoothness penalty, but suggests that a thicker cortex in the red areas is indicative of AD, which is not physiologically plausible. When we explicitly model geometric differences, this evidence seems to disappear. This suggests that there is a non-negligible dependence structure between the predictors modeling geometry and those modeling thickness. Differences that seemed to be related to cortical thickness in the model without the geometric component are now captured by the term that models cortical geometric variations. Furthermore, when we model intersubject geometric differences the entorhinal cortex atrophy in the medial temporal lobe is identified as the strongest predictor of AD. This is consistent with pathological findings and staging of early AD (Braak et al. 2006).

information due to misregistration. While incorporating geometric information into the model may lead to only minor improvements in classification performance, as shown in Figure 4, the estimated discriminant direction can be significantly different between the two models. Although the ground truth is unknown, the estimated discriminant direction when geometric information is included is more consistent with the neurodegenerative nature of the disease, as explained in the next section.

#### 6.6 Discussion

The results in Figure 3 identify the typical AD thickness signature. Several studies that focus on identifying AD-vulnerable areas include the regions found in our analysis (see, e.g., Bondareff et al. 1989; Dickerson et al. 2009; Sabuncu et al. 2011). However, there is some variability in the estimated regions. For instance, Sabuncu et al. (2011) used a dynamic model and found strongest changes in the inferior parietal regions and the posterior cingulate. It should be noted that these studies typically consist of massive univariate analyses between the cortical thickness at each vertex, or each parcel, and the diagnostic label. They are therefore taking a feature-centric perspective on the problem. It is not clear how these findings would generalize to out-of-sample data (Li and Tong 2020).

To demonstrate the importance of modeling cortical geometry, we compare our results to those obtained by fitting a ridge logistic regression model and the proposed model in equation (10), i.e., by discarding inter-subject geometric differences. We compare these estimates in Figure 4. What we observe is that ridge logistic regression yields estimated discriminant directions that are more difficult to interpret, due to the high spatial variability. Except for the entorhinal cortex, the functional model in equation (10) is able to capture the main areas where cortical thinning is associated with AD.

However, this model also suggests that a thicker cortex in certain regions (dark red) is associated with AD, contradicting the neurodegenerative nature of AD. Interestingly, introducing the geometric component in the model reduces such effects. This may also be caused by the geometric component now capturing systematic misregistration. In order to verify such a hypothesis, further validation of the estimated geometric component is required in controlled settings where registration is more reliable, e.g., in the longitudinal setting.

### 7 Conclusions

We introduce a framework for the discriminant analysis of functional data supported on random manifold domains, i.e., FoSs. To this aim, we adopt linear representations of these objects that are bivariate functional data belonging to linear spaces. We then define a functional linear classification model on the parametrizing space. Thanks to a penalized least-squares formulation, the proposed model is able to estimate the most discriminant direction in the data without requiring the explicit computation of the covariance function of the predictors or low-rank approximations thereof. This allows us to reduce the memory requirements by five orders of magnitude and ultimately be able to run our model on a standard workstation. The complexity of the solution is controlled by means of differential penalties that are aware of the geometry of the domain where the functional data are supported.

We apply the proposed model to the analysis of modern multi-modal neuroimaging data. Specifically, we estimate interpretable discriminant directions that are able to leverage both geometric and thickness features of the cerebral cortex to identify subjects with AD. Our results are consistent with those in the neuroscience literature.

The model proposed can be applied to several imaging settings that lead to FoSs representations, such as musculoskeletal imaging (Gee et al. 2018) or cardiac imaging (Biffi et al. 2018). It is also important to highlight that the proposed model is not a mere generalization of existing models for functional data supported on one-dimensional domains to multidimensional domains. We believe that its application to one-dimensional functional data, where the bivariate representation is given by the registered functions and associated registration maps, leads to a novel classification approach in this simplified setting.

## Acknowledgments

Data used in the preparation of this article were obtained from two sources: the Alzheimer's Disease Neuroimaging Initiative (ADNI) and the Parkinson's Progression Markers Initiative (PPMI). 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, Alzheimer's Association; Alzheimer's 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). The grantee organization is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer's 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.

PPMI — a public-private partnership — is funded by the Michael J. Fox Foundation for Parkinson's Research and funding partners. The complete list of PPMI funding partners can be found at www.ppmi-info.org.

ADNI data are available to the scientific community thorough the LONI Image and Data Archive at http://adni.loni.usc.edu/data-samples/accessdata. PPMI data can be accessed through https://www.ppmi-info.org/access-data-specimens/download-data.

# **Appendices**

### A Linear functional representation model

### A.1 Diffeomorphic deformation operator

The diffeomorphic operator  $\varphi$  can be constructed as follows. Let  $\{v_t \in \mathcal{V}(\mathbb{R}^3) : t \in [0,1]\}$  be a timevariant vector field such that  $\int_0^1 ||v_t||^2_{\mathcal{V}(\mathbb{R}^3)} dt < \infty$ . Then, the solution  $\phi_v : [0,1] \times \mathbb{R}^3 \to \mathbb{R}^3$ , at time t = 1, to the Ordinary Differential Equation (ODE)

$$\begin{cases} \frac{\partial \phi_v}{\partial t}(t, x) = v_t \circ \phi_v(t, x) & t \in [0, 1], x \in \mathbb{R}^3, \\ \phi_v(0, x) = x & x \in \mathbb{R}^3, \end{cases}$$
(22)

is a diffeomorphic deformation of  $\mathbb{R}^3$  (see, e.g., Younes 2019). We then model  $\{v_t : t \in [0,1]\}$  as a minimizer of the quantity  $\int_0^1 ||v_t||_{\mathcal{V}}^2 dt$ , for a given initial vector field  $v_0 \in \mathcal{V}(\mathbb{R}^3)$  (Miller et al. 2006). Finally, the diffeomorphic operator is defined as  $\varphi_{v_0}(x) = \varphi_v(1,x)$ , where  $v_0 \in \mathcal{V}(\mathbb{R}^3)$  is the initial vector field generating  $\{v_t : t \in [0,1]\}$ , and  $\phi_v$  is the solution to the ODE in equation (22) for the computed  $\{v_t : t \in [0,1]\}$ .

#### A.2 Computation

In practice, each surface  $\mathcal{M}_i$  has a computational representation  $\mathcal{M}_i^{\mathcal{T}}$  that is a triangle mesh with s vertices  $\xi_1^i, \ldots, \xi_s^i$  in correspondence across the n subjects. We use these vertices to perform Procrustes analysis (Dryden and Mardia 2016) and remove translation, size, and rigid rotations from the surfaces  $\mathcal{M}_i^{\mathcal{T}}$ . In addition, Procrustes analysis yields a template  $\mathcal{M}^{\mathcal{T}}$  with vertices  $\xi_1, \ldots, \xi_s$ .

The representation functions  $\{v_i : v_i \in \mathcal{V}(\mathbb{R}^3)\}$ , associated with the surfaces  $\{\mathcal{M}_i^{\mathcal{T}}\}$ , are then computed by solving the minimization problem

$$v_{i} = \arg\min_{v \in \mathcal{V}(\mathbb{R}^{3})} \sum_{l=1}^{s} \|\varphi_{v}(\xi_{l}) - \xi_{l}^{i}\|_{\mathbb{R}^{3}}^{2} + \lambda \|v\|_{\mathcal{V}(\mathbb{R}^{3})}^{2}, \qquad i = 1, \dots, n,$$
(23)

where the least-squares term ensures that the deformed template  $\varphi_{v_i}(\mathcal{M}^T)$  is a close approximation of  $\mathcal{M}_i^T$ . The term  $\|v\|_{\mathcal{V}(\mathbb{R}^3)}^2$  is a regularizing term that encourages the solution to achieve such a close approximation with a 'minimal' deformation. The constant  $\lambda$  is selected by inspecting the solutions on a small subset of the full cohort. To perform the actual computations, we use the MATLAB implementation fshapetk (Charlier et al. 2015; Charlier et al. 2017). Note that if the vertices were not in correspondence, that is, the surfaces had not been registered beforehand, the proposed framework would still be applicable by replacing the least-squares term in equation (23) with a more general shape similarity measure  $D(\cdot,\cdot)$ . An example of such a similarity measure is found in Vaillant and Glaunès (2005) and Vaillant et al. (2007), where the authors use the concept of currents, from geometric measure theory, to represent surfaces.

#### A.3 Selecting an appropriate kernel

The main requirement for choosing the kernel  $K_{\mathbb{R}^3}$  is that the associated space  $\mathcal{V}(\mathbb{R}^3)$  is an admissible space (Younes 2019). Therefore, there must exist a positive constant M such that for all  $v \in \mathcal{V}(\mathbb{R}^3)$ , the following inequality holds:

$$||v||_{1,\infty} \le M||v||_{\mathcal{V}(\mathbb{R}^3)},$$

where  $\|\cdot\|_{1,\infty}$  is the canonical norm of the space  $C^1(\mathbb{R}^3)$ . This condition guarantees that  $\varphi_v$  is diffeomorphic for any  $v \in \mathcal{V}(\mathbb{R}^3)$ .

Nonetheless, within the set of admissible spaces, the specific choice of the kernel can significantly influence the quality of the estimated representations. We have found the approach proposed in Bruveris

et al. (2012) to be effective. This uses a mixture of isotropic Gaussian kernels with different variances, which intuitively allows for a multi-scale representation of the surfaces. In our application, we use a mixture of six Gaussian kernels with variance parameters set to  $(\sigma_1^2, \ldots, \sigma_6^2) = (64, 16, 4, 1, 0.25, 0.01)$ . Our choice was informed by visual inspection of the differences between the estimated  $\varphi_{v_i}(\mathcal{M}^T)$  and  $\mathcal{M}_i^T$ .

#### A.4 Template estimation

```
Algorithm 1: Algorithm for template estimation.

Data: The surfaces \mathcal{M}_1, \ldots, \mathcal{M}_n and an initial guess for the template \hat{\mathcal{M}}^{\{1\}} = \mathcal{M}_1

Result: \hat{\mathcal{M}} = \hat{\mathcal{M}}^{\{N_{\text{iter}}\}}

for iter = 1, ..., N_{\text{iter}} do

| for i = 1, \ldots, n do

| \hat{v}_i^{\{\text{iter}\}} \leftarrow \arg\min_{v \in \mathcal{V}(\mathbb{R}^3)} D\left(\varphi_v(\mathcal{M}^{\{\text{iter}\}}), \mathcal{M}_i\right) + \lambda ||v||_{\mathcal{V}(\mathbb{R}^3)}^2

end

| \hat{\mu}^{\{\text{iter}\}} = \frac{1}{n} \sum_{i=1}^n \hat{v}_i^{\{\text{iter}\}}
| \hat{\mathcal{M}}^{\{\text{iter}\}} = \varphi_{\hat{\mu}^{\{\text{iter}\}}}(\mathcal{M})

end
```

In this section, we introduce an algorithm designed to estimate a template leveraging the (formal) Riemannian structure of the manifold of diffeomorphisms adopted to model random manifold domains. One approach is to define the template so that the average of the linear representations  $\{v_i\}$ , located on the tangent space at the identity map, is zero. The details of this iterative centroid approach are outlined in Algorithm 1. For an overview of alternative approaches see Cury et al. (2014).

Despite leveraging GPU acceleration for computing RKHS norms and associated gradients, the process of computing  $v_i$  for each subject still takes about 40 minutes in our application. This makes the process of estimating the template computationally prohibitive given the necessity of multiple iterations. Therefore, we have chosen to use a fixed template in our final application.

#### **B** Simulations

In this section, we conduct simulations to assess the finite sample classification performance of the model proposed when compared to the models introduced in Section 6.5. Here we focus on the functional univariate setting described in Section 3.

We use a triangle mesh  $\mathcal{M}_{\mathcal{T}}$  with 642 nodes that is an approximation of a brainstem. On this triangulated surface, we generate the orthonormal functions  $\{v_l : l = 1, 2, ..., 40\}$  consisting of 40 eigenfunctions of the Laplace-Beltrami operator computed on  $\mathcal{M}_{\mathcal{T}}$ . Then, we generate two sets of smooth functional data supported on  $\mathcal{M}_{\mathcal{T}}$ , with identical within-group covariance structures, as follows:

Group 1: 
$$x_{i1} = w_{i1}v_1 + w_{i2}v_2 + \ldots + w_{i40}v_{40} \quad i = 1, \ldots, n,$$
  
Group 2:  $x_{i2} = \alpha\mu + u_{i1}v_1 + u_{i2}v_2 + \ldots + u_{i40}v_{40} \quad i = 1, \ldots, n,$  (24)

where  $u_{ij}$  and  $w_{ij}$  are zero-mean independent random variables that represent the scores and are distributed according to a normal distribution with variance  $\sigma_j^2$  decreasing in j. The function  $\mu$  is the groups' difference, chosen to be a fixed linear combination of the eigenfunctions  $\{v_l : l = 1, 2, \ldots, 40\}$ , and its magnitude is controlled by a parameter  $\alpha > 0$  defining the 'difficulty' of the classification problem, that is, the signal-to-noise ratio. We then identify three regimes, with  $\alpha \in \{0.2, 0.4, 0.6\}$ . Given a new function  $x^*$ , the aim is to recover the group this observation belongs to.

We compare the proposed FLDA method against the models introduced in Section 6.5, e.g., FPCA followed by LDA, Lasso and Ridge logistic regression, random forests, support vector machines, and

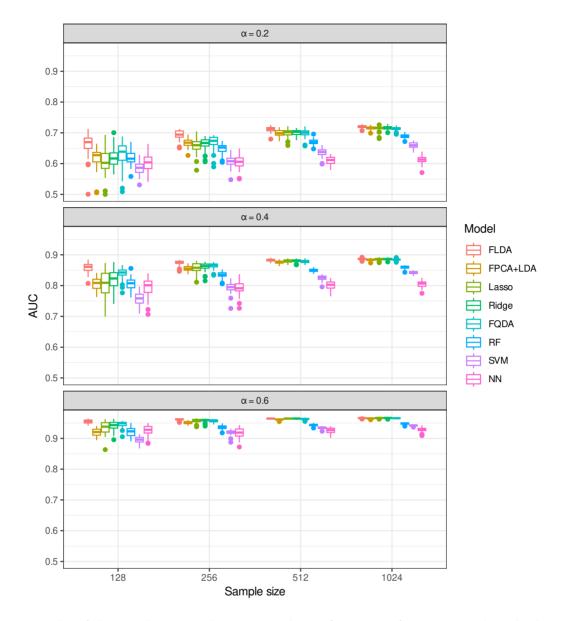


Figure 5: Results of the simulation study to assess the performance of our proposed method, under the assumption of homogeneous covariances, for various sample sizes (n=128,256,512,1024) and signal-to-noise ratios ( $\alpha=0.2,0.4,0.6$ ), where  $\alpha$  reflects the strength of the discriminant signal. Prediction accuracy is measured using AUC and the simulations were repeated 50 times for each setting.

fully-connected feed-forward neural networks. For standard multivariate models, we use the values of  $x_i$  at the vertices of  $\mathcal{M}_{\mathcal{T}}$  to construct the data matrix. We evaluate the performance of each method for different sample sizes of the training data: n=128,256,512, and 1024. For every n, besides the training set, we generate a validation set of size n and a test set including 20K samples, and repeat the experiment 50 times. To select the hyperparameters of the models, we employ a validation set approach. We summarize the classification performances on the test set in Figure 5.

Figure 5 shows that by increasing  $\alpha$ , or the sample size n, the performances of all models tend to improve. This is expected, given that a larger  $\alpha$  makes the classification task easier. In addition, a larger sample size allows for a more accurate estimation of the unknown model parameters. In the setting of equal within-class covariance structures, it is well known that the best classifier is linear, so it is not surprising that some of the nonlinear models, e.g., RF, SVM, and NN, show worse performances. In particular, FLDA appears to perform better than the other methods. For large sample sizes, the influence of the regularization terms becomes negligible and the difference in

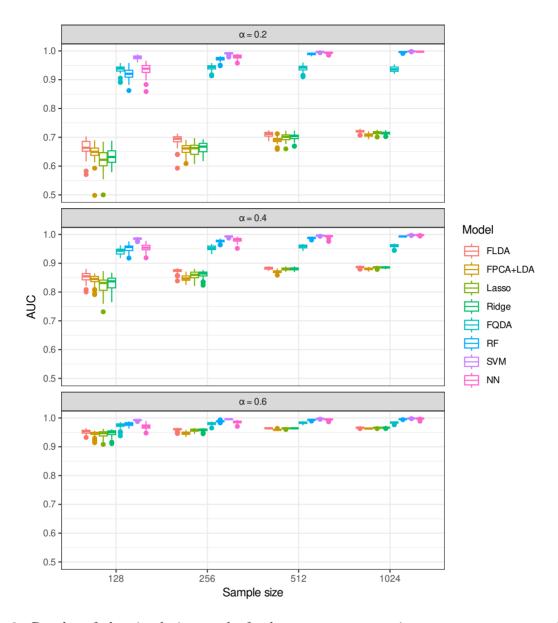


Figure 6: Results of the simulation study for heterogeneous covariance structures across different sample sizes (n=128,256,512,1024) and signal-to-noise ratios ( $\alpha=0.2,0.4,0.6$ ), where  $\alpha$  reflects the strength of the discriminant signal. The prediction accuracy was evaluated through AUC and the simulations were repeated 50 times for each setting.

performance starts vanishing.

Next, we explore the performance of the different methods in the setting where the two groups have different within-group covariance structures. Specifically, we generate the data as in equation (24) with  $u_{ij}$  zero-mean independent random variables distributed according to a normal with variance  $\sigma_j^2$ , decreasing in j, and  $w_{ij}$  zero-mean independent random variables distributed according to a normal with variance  $\sigma_{40-j}^2$ , and therefore, increasing in j. The results of the simulations are shown in Figure 6.

In this setting, it is well known that the best classifier is not linear. Hence, it is not surprising that the linear methods tend to perform worse than the nonlinear ones and that this difference does not vanish by simply increasing the sample size. Moreover, a larger  $\alpha$  leads to better performance across all the tested methods. Specifically, the proposed FQDA model outperforms all the linear models, but other non-linear approaches, such as SVM, perform even better. Importantly, the proposed FLDA model performs best among the linear models even when it is misspecified.

Classification performance is not the only relevant metric. In the application motivating this

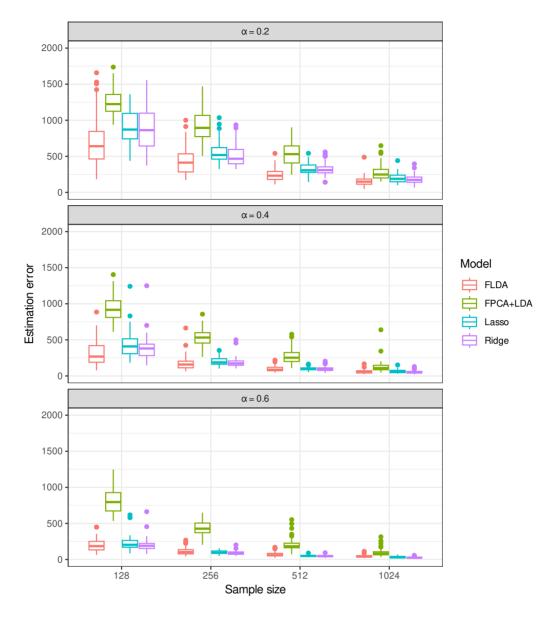


Figure 7: Results of the simulation study to compare the performance of the different linear methods considered, using homogeneous covariances, for various sample sizes (n = 128, 256, 512, 1024) and signal-to-noise ratios ( $\alpha = 0.2, 0.4, 0.6$ ). Here, we measure the performance using the estimation error  $\|\hat{\beta} - \beta^0\|_{\mathcal{L}^2(\mathcal{M})}^2$ , with  $\hat{\beta}$  an appropriately normalized version of the estimate of the true functional parameter  $\beta^0$ .

work, we are particularly interested in accurately estimating the classification rule, which for linear models is fully described by a parameter  $\hat{\beta}$ . Therefore, in the setting of homogeneous covariances, we compare the performance of the different models with respect to the metric  $\|\hat{\beta} - \beta^0\|_{\mathcal{L}^2(\mathcal{M})}^2$ , with  $\hat{\beta}$  an appropriately normalized version of the estimate of the true functional parameter  $\beta^0$ . Note that the estimates from a linear discriminant analysis and a logistic regression model can be compared, as they assume the same model but estimate their parameters differently (Hastie et al. 2009). For the standard multivariate methods,  $\hat{\beta}$  is constructed by interpolating a piecewise linear function to its estimated discrete counterpart. The results, shown in Figure 7, indicate that the proposed FLDA model does not only yield more accurate predictions but also more accurate estimates of the underlying functional parameters.

### C Proofs

Recall that we denote by C the covariance function of the functional predictor and by  $\hat{C}$  its empirical counterpart. Moreover,  $K_{\mathcal{M}}$  denotes the evaluation kernel of the Sobolev space  $\mathcal{W}^2(\mathcal{M})$  endowed with the norm  $\left(\|\Delta_{\mathcal{M}}\beta\|_{\mathcal{L}^2(\mathcal{M})}^2 + \varepsilon\|\beta\|_{\mathcal{L}^2(\mathcal{M})}^2\right)^{1/2}$ . The sandwich operator T is defined as  $T = L_{K_{\mathcal{M}}}^{\frac{1}{2}} L_C L_{K_{\mathcal{M}}}^{\frac{1}{2}}$ . To simplify the notation, we drop the subscripts from  $\|\cdot\|_{\mathcal{L}^2(\mathcal{M})}$  and  $\langle\cdot,\cdot\rangle_{\mathcal{L}^2(\mathcal{M})}$ . We define  $\|A\|_{\mathrm{op}} = \sup_{f:\|f\|=1} \|Af\|$  to be the operator norm of a linear operator  $A: \mathcal{L}^2(\mathcal{M}) \to \mathcal{L}^2(\mathcal{M})$ .

Thanks to the fact that  $L_{K_{\mathcal{M}}}^{\frac{1}{2}}(\mathcal{L}^2(\mathcal{M})) = \mathcal{W}^2(\mathcal{M})$  (Cucker and Smale 2002), it is clear that we can reformulate the problem in equation (10) as

$$\operatorname{minimize}_{f} \frac{1}{n} \sum_{i=1}^{n} \left( y_{i} - \langle x_{i}, L_{K_{\mathcal{M}}}^{\frac{1}{2}} f \rangle \right)^{2} + \lambda \|f\|^{2}, \tag{25}$$

whose solution  $\hat{f}_{\lambda}$  is given by

$$\hat{f}_{\lambda} = (T_n + \lambda I)^{-1} \frac{1}{n} \sum_{i=1}^{n} y_i L_{K_{\mathcal{M}}}^{\frac{1}{2}} x_i,$$
(26)

where  $T_n = L_{K_{\mathcal{M}}}^{\frac{1}{2}} L_{\hat{C}} L_{K_{\mathcal{M}}}^{\frac{1}{2}}$ . Moreover, the out-of-sample risk  $\mathbb{E}^* \left[ \langle X^*, \beta^0 - \hat{\beta} \rangle \right]^2$  can be easily rewritten as  $\left\| T^{\frac{1}{2}} (\hat{f}_{\lambda} - f_0) \right\|^2$ , with  $f_0 \in \mathcal{L}^2(\mathcal{M})$  such that  $f_0 = T^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} (\mu_2 - \mu_1)$ . Observe that

$$\begin{split} &T^{\frac{1}{2}}(\hat{f}_{\lambda} - f_{0}) \\ &= T^{\frac{1}{2}} \left[ (T_{n} + \lambda I)^{-1} \left( \frac{1}{n} \sum_{i=1}^{n} y_{i} L_{K_{\mathcal{M}}}^{\frac{1}{2}} x_{i} \right) - f_{0} \right] \\ &= T^{\frac{1}{2}} \left[ (T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - (\mu_{2} - \mu_{1}) \right) + (T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} (\mu_{2} - \mu_{1}) - T^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} (\mu_{2} - \mu_{1}) \right] \end{split}$$

Let  $\hat{d} = \frac{1}{n} \sum_{i=1}^{n} y_i x_i$  and  $d = \mu_2 - \mu_1$ , and notice that  $\mathbb{E}[\hat{d}] = d + o(1)$ . Then, we have

$$\begin{aligned}
& \left\| T^{\frac{1}{2}}(\hat{f}_{\lambda} - f_{0}) \right\| \\
& \leq \left\| T^{\frac{1}{2}}(T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\| + \left\| T^{\frac{1}{2}} \left[ (T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d - T^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d \right] \right\| \\
& = I_{1} + I_{2}
\end{aligned} (27)$$

We first derive a bound for the term  $I_1$  in equation (27), i.e., the variance term, and then proceed with bounding the term  $I_2$ , i.e., the bias term. To accomplish this, we will also use Theorems C.2-C.3, stated in Section C.4.

#### C.1 Variance

Simple calculations show that

$$\begin{split} I_{1} &= \left\| T^{\frac{1}{2}} (T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\| \\ &= \left\| T^{\frac{1}{2}} (T_{n} + \lambda I)^{-\frac{1}{2}} (T_{n} + \lambda I)^{-\frac{1}{2}} (T + \lambda I)^{\frac{1}{2}} (T + \lambda I)^{-\frac{1}{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\| \\ &\leq \left\| T^{\frac{1}{2}} (T_{n} + \lambda I)^{-\frac{1}{2}} \right\|_{\text{op}} \left\| (T_{n} + \lambda I)^{-\frac{1}{2}} (T + \lambda I)^{\frac{1}{2}} \right\|_{\text{op}} \left\| (T + \lambda I)^{-\frac{1}{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\| \\ &\leq \left\| (T + \lambda I)^{\frac{1}{2}} (T_{n} + \lambda I)^{-\frac{1}{2}} \right\|_{\text{op}} \left\| (T_{n} + \lambda I)^{-\frac{1}{2}} (T + \lambda I)^{\frac{1}{2}} \right\|_{\text{op}} \left\| (T + \lambda I)^{-\frac{1}{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\| \\ &\leq \left\| (T + \lambda I) (T_{n} + \lambda I)^{-1} \right\|_{\text{op}} \left\| (T + \lambda I)^{-\frac{1}{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\|, \end{split}$$

where in the last inequality we have used  $||A^{\gamma}B^{\gamma}||_{\text{op}} \leq ||AB||_{\text{op}}^{\gamma}$ , for any  $0 < \gamma < 1$  (Blanchard and Krämer 2010). We can then bound the first term thanks to Theorem C.3. We therefore turn to the second term.

First observe that

$$\mathbb{E}\left[\hat{d}\right] = d + o(1), \qquad \mathbb{E}\left[(\hat{d} - d) \otimes (\hat{d} - d)\right] = \frac{1}{n}\left(C - \pi_1\mu_1 \otimes \mu_1 - \pi_2\mu_2 \otimes \mu_2\right) + o(1),$$

and, therefore,

$$\begin{split} & \langle L_{K_{\mathcal{M}}}^{\frac{1}{2}}(\hat{d} - d), \eta_{k} \rangle^{2} \\ &= \langle L_{K_{\mathcal{M}}}^{\frac{1}{2}} L_{\mathbb{E}[(\hat{d} - d) \otimes (\hat{d} - d)]} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \eta_{k}, \eta_{k} \rangle \\ &= \frac{1}{n} \langle T \eta_{k}, \eta_{k} \rangle - \frac{\pi_{1}}{n} \langle L_{K_{\mathcal{M}}}^{\frac{1}{2}} L_{\mu_{1} \otimes \mu_{1}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \eta_{k}, \eta_{k} \rangle - \frac{\pi_{2}}{n} \langle L_{K_{\mathcal{M}}}^{\frac{1}{2}} L_{\mu_{2} \otimes \mu_{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \eta_{k}, \eta_{k} \rangle \\ &= \frac{1}{n} \left( \tau_{k} - \pi_{1} \| L_{K_{\mathcal{M}}}^{\frac{1}{2}} \mu_{1} \|^{2} - \pi_{2} \| L_{K_{\mathcal{M}}}^{\frac{1}{2}} \mu_{2} \|^{2} \right), \end{split}$$

with  $\{\tau_k\}$  and  $\{\eta_k\}$  denoting the eigenvalues and eigenfunctions of T, respectively. As in Gaynanova and Kolar (2015), we ignore the bias term o(1).

We then have

$$\mathbb{E} \left\| (T + \lambda I)^{-\frac{1}{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\|^{2}$$

$$= \mathbb{E} \left\| \sum_{k} (T + \lambda I)^{-\frac{1}{2}} \eta_{k} \langle L_{K_{\mathcal{M}}}^{\frac{1}{2}} (\hat{d} - d), \eta_{k} \rangle \right\|^{2}$$

$$\leq \frac{1}{n} \sum_{k} \left( \frac{\tau_{k}}{\lambda + \tau_{k}} - \pi_{1} \frac{\|L_{K_{\mathcal{M}}}^{\frac{1}{2}} \mu_{1}\|^{2}}{\lambda + \tau_{k}} - \pi_{2} \frac{\|L_{K_{\mathcal{M}}}^{\frac{1}{2}} \mu_{2}\|^{2}}{\lambda + \tau_{k}} \right)$$

$$\leq \frac{1}{n} D(\lambda).$$

Moreover, by appealing to the Markov inequality, we have that with confidence at least  $1 - \delta/2$ 

$$\left\| (T + \lambda I)^{-\frac{1}{2}} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\| \le \frac{2}{\delta} \sqrt{\frac{D(\lambda)}{n}} \le \frac{1}{\kappa \delta} B_{n,\lambda}, \tag{28}$$

where  $B_{n,\lambda} = \frac{2\kappa}{\sqrt{n}} \left( \frac{\kappa}{\sqrt{n\lambda}} + \sqrt{D(\lambda)} \right)$  and  $\kappa$  is defined such that  $\kappa^2 = \operatorname{ess\,sup} \|L_{K_{\mathcal{M}}}^{1/2} X\|^2$ .

Thanks to the inequality in equation (28) and Theorem C.3, with probability at least  $1 - \delta$ , we have

$$I_{1} = \left\| T^{\frac{1}{2}} (T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \frac{1}{n} \sum_{i=1}^{n} y_{i} x_{i} - d \right) \right\| \leq \left( \frac{B_{n,\lambda} \log(2/\delta)}{\lambda} + 1 \right)^{2} \left( \frac{1}{2\kappa \delta} B_{n,\lambda} \right). \tag{29}$$

We now turn to the bias term  $I_2$  in equation (27).

#### C.2 Bias

By simple calculations, we have

$$\begin{split} I_{2} &= \left\| T^{\frac{1}{2}} \left[ (T_{n} + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d - T^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d \right] \right\| \\ &= \left\| T^{\frac{1}{2}} (T_{n} + \lambda I)^{-1} [T - T_{n} + \lambda I] \right\|_{\text{op}} \|f_{0}\| \\ &\leq \left\| (T + \lambda I)^{\frac{1}{2}} (T_{n} + \lambda I)^{-\frac{1}{2}} \right\|_{\text{op}} \left\| (T_{n} + \lambda I)^{-\frac{1}{2}} (T + \lambda I)^{\frac{1}{2}} \right\|_{\text{op}} \left\| (T + \lambda I)^{-\frac{1}{2}} [T - T_{n} + \lambda I] \right\|_{\text{op}} \|f_{0}\| \\ &\leq \left\| (T + \lambda I) (T_{n} + \lambda I)^{-1} \right\|_{\text{op}} \left\| (T + \lambda I)^{-\frac{1}{2}} [T - T_{n}] \right\|_{\text{op}} \|f_{0}\| \\ &\leq \left\| (T + \lambda I) (T_{n} + \lambda I)^{-1} \right\|_{\text{op}} \left\| (T + \lambda I)^{-\frac{1}{2}} [T - T_{n}] \right\|_{\text{op}} \|f_{0}\| \\ &+ \left\| (T + \lambda I) (T_{n} + \lambda I)^{-1} \right\|_{\text{op}} \lambda \left\| (T + \lambda I)^{-\frac{1}{2}} \right\|_{\text{op}} \|f_{0}\|. \end{split}$$

Moreover, by using Theorems C.2-C.3, and the inequality

$$\lambda \left\| (T + \lambda I)^{-\frac{1}{2}} \right\|_{\text{op}} \le \frac{\lambda}{\sqrt{\lambda}} = \sqrt{\lambda},\tag{30}$$

we have that with probability at least  $1 - \delta$ 

$$\left\| T^{\frac{1}{2}} \left[ (T_n + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d - T^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d \right] \right\| \le \left( \frac{B_{n,\lambda} \log(4/\delta)}{\sqrt{\lambda}} + 1 \right)^2 (B_{n,\lambda} \log(4/\delta) + \sqrt{\lambda}) \|f_0\|. \tag{31}$$

#### C.3 Final rates

Using the variance and bias bounds in equations (29) and (31), we get the following bound for the out-of-sample risk:

$$\left\| T^{\frac{1}{2}}(\hat{f}_{\lambda} - f_0) \right\|^2$$
 (32)

$$\leq 2 \left\| T^{\frac{1}{2}} (T_n + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} \left( \hat{d} - d \right) \right\|^2 + 2 \left\| T^{\frac{1}{2}} \left[ (T_n + \lambda I)^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d - T^{-1} L_{K_{\mathcal{M}}}^{\frac{1}{2}} d \right] \right\|^2 \tag{33}$$

$$\leq 2 \left( \frac{B_{n,\lambda} \log(2/\delta)}{\sqrt{\lambda}} + 1 \right)^4 \left( \frac{1}{2\kappa\delta} B_{n,\lambda} \right)^2 \tag{34}$$

$$+2\left(\frac{B_{n,\lambda}\log(4/\delta)}{\sqrt{\lambda}}+1\right)^4\left(B_{n,\lambda}\log(4/\delta)+\sqrt{\lambda}\right)^2\|f_0\|^2\tag{35}$$

$$\leq 2 \frac{\lambda}{\delta^2} \left( \frac{B_{n,\lambda} \log(2/\delta)}{\sqrt{\lambda}} + 1 \right)^4 \left( \frac{1}{2\kappa} \frac{B_{n,\lambda}}{\sqrt{\lambda}} \right)^2 \tag{36}$$

$$+2\lambda \left(\frac{B_{n,\lambda}\log(4/\delta)}{\sqrt{\lambda}}+1\right)^4 \left(\frac{B_{n,\lambda}}{\sqrt{\lambda}}\log(4/\delta)+1\right)^2 \|f_0\|^2$$
(37)

$$\leq C \frac{(\log(4/\delta))^6}{\delta^2} n^{-\frac{1}{1+\theta}},$$
(38)

where in the last inequality we have chosen  $\lambda = n^{-\frac{1}{1+\theta}}$  and used the inequality  $B_{n,\lambda} \leq 2\kappa(\kappa + \sqrt{c})\sqrt{\lambda}$ . The constant c here is from Assumption 3.3.

This implies the result stated in Theorem 3.1.

#### C.4 Auxiliary results

**Theorem C.1.** Let  $\mathcal{H}$  be a Hilbert space endowed with a norm  $\|\cdot\|_{\mathcal{H}}$  and let X be a random variable taking values in  $\mathcal{H}$ . Let  $\xi_1, \ldots, \xi_n$  be a sequence of n independent copies of X. Assume that  $\|\xi\|_{\mathcal{H}} \leq M$  (a.s.), then for  $0 < \delta < 1$ 

$$\left\| \frac{1}{n} \sum_{i=1}^{n} \left( \xi_i - \mathbb{E}[\xi] \right) \right\|_{\mathcal{U}}^2 \le \frac{2M \log(2/\delta)}{n} + \sqrt{\frac{2\mathbb{E}\left[ \|\xi\|_{\mathcal{H}}^2 \right] \log(2/\delta)}{n}}$$

with probability at least  $1 - \delta$ .

*Proof.* A proof can be found in Pinelis (2007).

Using Theorem C.1, it can be shown that the following two theorems hold.

**Theorem C.2.** Under Assumption 3.1, for any  $0 < \delta < 1$ , the inequality

$$\|(T + \lambda I)^{-\frac{1}{2}}(T - T_n)\|_{\text{op}} \le B_{n,\lambda} \log(2/\delta)$$
 (39)

holds with confidence at least  $1 - \delta$ .

*Proof.* A proof can be found in Tong and Ng (2018).

**Theorem C.3.** Under Assumption 3.1, for any  $0 < \delta < 1$  with confidence at least  $1 - \delta$ ,

$$\|(T+\lambda I)(T_n+\lambda I)^{-1}\|_{\text{op}} \le \left(\frac{B_{n,\lambda}\log(2/\delta)}{\sqrt{\lambda}}+1\right)^2.$$
(40)

Moreover, the confidence set is the same as the one in Theorem C.2.

*Proof.* A proof can be found in Tong and Ng (2018).

#### C.5 Multivariate model

Thanks to the multivariate notation introduced in Section 4, we can reformulate the multivariate model in equation (15) as

$$\operatorname{minimize}_{\boldsymbol{f}} \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \langle \boldsymbol{x}_i, L_{\boldsymbol{K}}^{\frac{1}{2}} \boldsymbol{f} \rangle_{\mathcal{H}} \right)^2 + \lambda \|\boldsymbol{f}\|_{\mathcal{H}}^2, \tag{41}$$

with  $x_i := (v_i, x_i)$ . Therefore, the proof of Theorem 4.1 follows the same lines as that of Theorem 3.1 and is therefore omitted.

#### References

Arguillère, S., Miller, M. I., and Younes, L. (2016). "Diffeomorphic Surface Registration with Atrophy Constraints". In: *SIAM Journal on Imaging Sciences* 9.3, pp. 975–1003. ISSN: 19364954. DOI: 10.1137/15M104431X.

Berlinet, A. and Thomas-Agnan, C. (2004). Reproducing Kernel Hilbert Spaces in Probability and Statistics. Boston, MA: Springer US. ISBN: 978-1-4613-4792-7. DOI: 10.1007/978-1-4419-9096-9.

- Berrendero, J. R., Cuevas, A., and Torrecilla, J. L. (2018). "On the Use of Reproducing Kernel Hilbert Spaces in Functional Classification". In: *Journal of the American Statistical Association* 113.523, pp. 1210–1218. ISSN: 1537274X. DOI: 10.1080/01621459.2017.1320287.
- Biffi, C. et al. (2018). "Three-Dimensional Cardiovascular Imaging-Genetics: A Mass Univariate Framework". In: *Bioinformatics* 34.1, pp. 97–103. ISSN: 14602059. DOI: 10.1093/bioinformatics/btx552.
- Blanchard, G. and Krämer, N. (2010). "Optimal Learning Rates for Kernel Conjugate Gradient Regression". In: Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010, NIPS 2010. ISBN: 978-1-61782-380-0. DOI: 10.1.1.231.648.
- Bondareff, W., Mountjoy, C. Q., Roth, M., and Hauser, D. L. (1989). "Neurofibrillary Degeneration and Neuronal Loss in Alzheimer's Disease". In: *Neurobiology of Aging* 10.6, pp. 709–715. ISSN: 01974580. DOI: 10.1016/0197-4580(89)90007-9.
- Braak, H., Alafuzoff, I., Arzberger, T., Kretzschmar, H., and Tredici, K. (2006). "Staging of Alzheimer Disease-Associated Neurofibrillary Pathology Using Paraffin Sections and Immunocytochemistry". In: *Acta Neuropathologica* 112.4, pp. 389–404. ISSN: 00016322. DOI: 10.1007/s00401-006-0127-z.
- Breiman, L. (2001). "Random Forests". In: *Machine Learning* 45.1, pp. 5–32. ISSN: 1573-0565. DOI: 10.1023/A:1010933404324.
- Brezis, H. (2011). Functional Analysis, Sobolev Spaces and Partial Differential Equations. New York, NY: Springer New York. ISBN: 978-0-387-70913-0. DOI: 10.1007/978-0-387-70914-7.
- Bruveris, M., Risser, L., and Vialard, F. X. (2012). "Mixture of Kernels and Iterated Semidirect Product of Diffeomorphisms Groups". In: *Multiscale Modeling and Simulation* 10.4, pp. 1344–1368. ISSN: 15403459. DOI: 10.1137/110846324.
- Cai, T. T. and Yuan, M. (2012). "Minimax and Adaptive Prediction for Functional Linear Regression". In: *Journal of the American Statistical Association* 107.499, pp. 1201–1216. ISSN: 0162-1459, 1537-274X. DOI: 10.1080/01621459.2012.716337.
- Charlier, B., Charon, N., and Trouvé, A. (2017). "The Fshape Framework for the Variability Analysis of Functional Shapes". In: *Foundations of Computational Mathematics* 17.2, pp. 287–357. ISSN: 16153383. DOI: 10.1007/s10208-015-9288-2.
- Charlier, B., Nardi, G., and Trouvé, A. (2015). "The Matching Problem between Functional Shapes via a BV-penalty Term: A \$\Gamma\$-Convergence Result". In: pp. 1–31.
- Chen, L. H. and Jiang, C. R. (2018). "Sensible Functional Linear Discriminant Analysis". In: *Computational Statistics and Data Analysis* 126, pp. 39–52. ISSN: 01679473. DOI: 10.1016/j.csda. 2018.04.005.
- Chung, M. K., Hanson, J. L., and Pollak, S. D. (2021). "Statistical Analysis on Brain Surfaces". In: *Handbook of Neuroimaging Data Analysis*, pp. 273–302. DOI: 10.1201/9781315373652-20.
- Chung, M. K., Hartley, R., Dalton, K. M., and Davidson, R. J. (2008). "Encoding Cortical Surface by Spherical Harmonics". In: *Statistica Sinica* 18.4, pp. 1269–1291. ISSN: 10170405.
- Chung, M. K., Qiu, A., Seo, S., and Vorperian, H. K. (2015). "Unified Heat Kernel Regression for Diffusion, Kernel Smoothing and Wavelets on Manifolds and Its Application to Mandible Growth Modeling in CT Images". In: *Medical Image Analysis* 22.1, pp. 63–76. ISSN: 13618423. DOI: 10.1016/j.media.2015.02.003.
- Cortes, C. and Vapnik, V. (1995). "Support-Vector Networks". In: *Machine Learning* 20.3, pp. 273–297. ISSN: 0885-6125, 1573-0565. DOI: 10.1007/BF00994018.
- Cucker, F. and Smale, S. (2002). "On the Mathematical Foundations of Learning". In: Bulletin of the American Mathematical Society 39.1, pp. 1–49. ISSN: 02730979. DOI: 10.1090/S0273-0979-01-00923-5.
- Cury, C., Glaunès, J. A., and Colliot, O. (2014). "Diffeomorphic Iterative Centroid Methods for Template Estimation on Large Datasets". In: *Geometric Theory of Information*. Ed. by F. Nielsen. Cham: Springer International Publishing, pp. 273–299. ISBN: 978-3-319-05316-5 978-3-319-05317-2. DOI: 10.1007/978-3-319-05317-2\_10.

- Dai, X. and Müller, H. G. (2018). "Principal Component Analysis for Functional Data on Riemannian Manifolds and Spheres". In: *Annals of Statistics* 46.6B, pp. 3334–3361. ISSN: 00905364. DOI: 10. 1214/17-AOS1660.
- Dai, X., Müller, H. G., and Yao, F. (2017). "Optimal Bayes Classifiers for Functional Data and Density Ratios". In: *Biometrika* 104.3, pp. 545–560. ISSN: 14643510. DOI: 10.1093/biomet/asx024.
- Dale, A. M., Fischl, B., and Sereno, M. I. (1999). "Cortical Surface-Based Analysis: I. Segmentation and Surface Reconstruction". In: *NeuroImage* 9.2, pp. 179–194. ISSN: 10538119. DOI: 10.1006/nimg.1998.0395.
- Delaigle, A. and Hall, P. (2012). "Achieving near Perfect Classification for Functional Data". In: *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 74.2, pp. 267–286. ISSN: 13697412. DOI: 10.1111/j.1467-9868.2011.01003.x.
- Dickerson, B. C. et al. (2009). "The Cortical Signature of Alzheimer's Disease: Regionally Specific Cortical Thinning Relates to Symptom Severity in Very Mild to Mild AD Dementia and Is Detectable in Asymptomatic Amyloid-Positive Individuals". In: Cerebral Cortex 19.3, pp. 497–510. ISSN: 10473211. DOI: 10.1093/cercor/bhn113.
- Dong, Q. et al. (2019). "Applying Surface-Based Hippocampal Morphometry to Study APOE-E4 Allele Dose Effects in Cognitively Unimpaired Subjects". In: *NeuroImage : Clinical* 22, p. 101744. ISSN: 2213-1582. DOI: 10.1016/j.nicl.2019.101744.
- Dong, Q. et al. (2020). "Applying Surface-Based Morphometry to Study Ventricular Abnormalities of Cognitively Unimpaired Subjects Prior to Clinically Significant Memory Decline". In: NeuroImage: Clinical 27, p. 102338. ISSN: 2213-1582. DOI: 10.1016/j.nicl.2020.102338.
- Dryden, I. L. and Mardia, K. V. (2016). Statistical Shape Analysis, with Applications in R: Second Edition. Chichester, UK: John Wiley & Sons, Ltd. ISBN: 978-1-119-07249-2. DOI: 10.1002/9781119072492.
- Dubey, P. and Müller, H. G. (2020). "Functional Models for Time-Varying Random Objects". In: *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 82.2, pp. 275–327. ISSN: 14679868. DOI: 10.1111/rssb.12337.
- Epifanio, I. and Ventura-Campos, N. (2014). "Hippocampal Shape Analysis in Alzheimer's Disease Using Functional Data Analysis". In: *Statistics in Medicine* 33.5, pp. 867–880. ISSN: 02776715. DOI: 10.1002/sim.5968.
- Fasshauer, G. E. and Ye, Q. (2013). Reproducing Kernels of Sobolev Spaces via a Green Kernel Approach with Differential Operators and Boundary Operators. Vol. 38. ISBN: 1-312-45173-4. DOI: 10.1007/s10444-011-9264-6.
- Feng, X., Li, T., Song, X., and Zhu, H. (2020). "Bayesian Scalar on Image Regression With Nonignorable Nonresponse". In: *Journal of the American Statistical Association* 115.532, pp. 1574–1597. ISSN: 1537274X. DOI: 10.1080/01621459.2019.1686391.
- Feragen, A., Lauze, F., and Hauberg, S. (2015). "Geodesic Exponential Kernels: When Curvature and Linearity Conflict". In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Vol. 07-12-June. IEEE, pp. 3032–3042. ISBN: 978-1-4673-6964-0. DOI: 10.1109/CVPR.2015.7298922.
- Ferrando, L., Ventura-Campos, N., and Epifanio, I. (2020). "Detecting and Visualizing Differences in Brain Structures with SPHARM and Functional Data Analysis". In: *NeuroImage* 222.August, p. 117209. ISSN: 10959572. DOI: 10.1016/j.neuroimage.2020.117209.
- Ferraty, F. and Vieu, P. (2006). *Nonparametric Functional Data Analysis*. Springer New York. ISBN: 978-0-387-30369-7. DOI: 10.1007/0-387-36620-2.
- Feydy, J., Charlier, B., Vialard, F. X., and Peyré, G. (2017). "Optimal Transport for Diffeomorphic Registration". In: Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). Ed. by M. Descoteaux, L. Maier-Hein, A. Franz, P. Jannin, D. Collins, and S. Duchesne. Vol. 10433 LNCS. Springer, Cham, pp. 291–299. ISBN: 978-3-319-66181-0. DOI: 10.1007/978-3-319-66182-7\_34.

- Fischl, B. and Dale, A. M. (2000). "Measuring the Thickness of the Human Cerebral Cortex from Magnetic Resonance Images". In: *Proceedings of the National Academy of Sciences of the United States of America* 97.20, pp. 11050–11055. ISSN: 00278424. DOI: 10.1073/pnas.200033797.
- Fischl, B., Sereno, M. I., and Dale, A. M. (1999). "Cortical Surface-Based Analysis: II. Inflation, Flattening, and a Surface-Based Coordinate System". In: *NeuroImage* 9.2, pp. 195–207. ISSN: 10538119. DOI: 10.1006/nimg.1998.0396.
- Fried, I. and Malkus, D. S. (1975). "Finite Element Mass Matrix Lumping by Numerical Integration with No Convergence Rate Loss". In: *International Journal of Solids and Structures* 11.4, pp. 461–466. ISSN: 00207683. DOI: 10.1016/0020-7683(75)90081-5.
- Gaynanova, I. (2020). "Prediction and Estimation Consistency of Sparse Multi-Class Penalized Optimal Scoring". In: *Bernoulli* 26.1, pp. 286–322. ISSN: 13507265. DOI: 10.3150/19-BEJ1126.
- Gaynanova, I. and Kolar, M. (2015). "Optimal Variable Selection in Multi-Group Sparse Discriminant Analysis". In: *Electronic Journal of Statistics* 9.2, pp. 2007–2034. ISSN: 19357524. DOI: 10.1214/15-EJS1064.
- Gaynanova, I. and Wang, T. (2019). "Sparse Quadratic Classification Rules via Linear Dimension Reduction". In: *Journal of Multivariate Analysis* 169, pp. 278–299. ISSN: 10957243. DOI: 10.1016/j.jmva.2018.09.011.
- Gee, A. H., Treece, G. M., and Poole, K. E. (2018). "How Does the Femoral Cortex Depend on Bone Shape? A Methodology for the Joint Analysis of Surface Texture and Shape". In: *Medical Image Analysis* 45, pp. 55–67. ISSN: 13618423. DOI: 10.1016/j.media.2018.01.001.
- Goldsmith, J., Huang, L., and Crainiceanu, C. M. (2014). "Smooth Scalar-on-Image Regression via Spatial Bayesian Variable Selection". In: *Journal of Computational and Graphical Statistics* 23.1, pp. 46–64. ISSN: 15372715. DOI: 10.1080/10618600.2012.743437.
- Grenander, U. and Miller, M. I. (1998). "Computational Anatomy: An Emerging Discipline". In: Quarterly of Applied Mathematics 56.4, pp. 617–694. ISSN: 0033-569X. DOI: 10.1090/qam/1668732.
- Happ, C. and Greven, S. (2018). "Multivariate Functional Principal Component Analysis for Data Observed on Different (Dimensional) Domains". In: *Journal of the American Statistical Association* 113.522, pp. 649–659. ISSN: 1537274X. DOI: 10.1080/01621459.2016.1273115.
- Hastie, T., Tibshirani, R., and Buja, A. (1994). "Flexible Discriminant Analysis by Optimal Scoring". In: *Journal of the American Statistical Association* 89.428, pp. 1255–1270. ISSN: 1537274X. DOI: 10.1080/01621459.1994.10476866.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). The Elements of Statistical Learning. New York, NY: Springer New York. ISBN: 978-0-387-84857-0. DOI: 10.1007/978-0-387-84858-7.
- Hazlett, H. C. et al. (2017). "Early Brain Development in Infants at High Risk for Autism Spectrum Disorder". In: *Nature* 542.7641, pp. 348–351. ISSN: 14764687. DOI: 10.1038/nature21369.
- Hinton, E., Rock, T., and Zienkiewicz, O. C. (1976). "A Note on Mass Lumping and Related Processes in the Finite Element Method". In: *Earthquake Engineering* \& Structural Dynamics 4.3, pp. 245—249. ISSN: 10969845. DOI: 10.1002/eqe.4290040305.
- Hoerl, A. E. and Kennard, R. W. (1970). "Ridge Regression: Biased Estimation for Nonorthogonal Problems". In: *Technometrics* 12.1, pp. 55–67. ISSN: 0040-1706, 1537-2723. DOI: 10.1080/00401706.1970.10488634.
- Horváth, L. and Kokoszka, P. (2012). Inference for Functional Data with Applications. Vol. 200. Springer Series in Statistics. New York, NY: Springer New York. ISBN: 978-1-4614-3654-6 978-1-4614-3655-3. DOI: 10.1007/978-1-4614-3655-3.
- Hsing, T. and Eubank, R. (2013). Theoretical Foundations of Functional Data Analysis, with an Introduction to Linear Operators. Chichester, UK: John Wiley & Sons, Ltd. ISBN: 978-1-118-76254-7. DOI: 10.1002/9781118762547.
- Im, K., Lee, J. M., Lyttelton, O., Kim, S. H., Evans, A. C., and Kim, S. I. (2008). "Brain Size and Cortical Structure in the Adult Human Brain". In: *Cerebral Cortex* 18.9, pp. 2181–2191. ISSN: 10473211. DOI: 10.1093/cercor/bhm244.

- Jack, C. R. et al. (2016). "A/T/N: An Unbiased Descriptive Classification Scheme for Alzheimer Disease Biomarkers". In: *Neurology* 87.5, pp. 539–547. ISSN: 1526632X. DOI: 10.1212/WNL.000000000002923.
- James, G. M. and Hastie, T. J. (2001). "Functional Linear Discriminant Analysis for Irregularly Sampled Curves". In: *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 63.3, pp. 533–550. ISSN: 13697412. DOI: 10.1111/1467-9868.00297.
- Jayasumana, S., Hartley, R., Salzmann, M., Li, H., and Harandi, M. (2015). "Kernel Methods on Riemannian Manifolds with Gaussian RBF Kernels". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37.12, pp. 2464–2477. ISSN: 01628828. DOI: 10.1109/TPAMI.2015.2414422.
- Jermyn, I. H., Kurtek, S., Klassen, E., and Srivastava, A. (2012). "Elastic Shape Matching of Parameterized Surfaces Using Square Root Normal Fields". In: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 7576 LNCS.PART 5, pp. 804–817. ISSN: 03029743. DOI: 10.1007/978-3-642-33715-4\_58.
- Jermyn, I. H., Kurtek, S., Laga, H., and Srivastava, A. (2017). "Elastic Shape Analysis of Three-Dimensional Objects". In: *Synthesis Lectures on Computer Vision* 7.3, pp. 1–185. ISSN: 2153-1056. DOI: 10.2200/s00785ed1v01v201707cov012.
- Jiang, C. R. and Wang, J. L. (2011). "Functional Single Index Models for Longitudinal Data". In: *Annals of Statistics* 39.1, pp. 362–388. ISSN: 00905364. DOI: 10.1214/10-AOS845.
- Kang, J., Reich, B. J., and Staicu, A. M. (2018). "Scalar-on-Image Regression via the Soft-Thresholded Gaussian Process". In: *Biometrika* 105.1, pp. 165–184. ISSN: 14643510. DOI: 10.1093/biomet/asx075.
- Kendall, D. G. (1984). "Shape Manifolds, Procrustean Metrics, and Complex Projective Spaces". In: Bulletin of the London Mathematical Society 16.2, pp. 81–121. ISSN: 14692120. DOI: 10.1112/blms/16.2.81.
- Kim, K. R., Dryden, I. L., Le, H., and Severn, K. E. (2021). "Smoothing Splines on Riemannian Manifolds, with Applications to 3D Shape Space". In: *Journal of the Royal Statistical Society.* Series B: Statistical Methodology 83.1, pp. 108–132. ISSN: 14679868. DOI: 10.1111/rssb.12402.
- Kraus, D. and Stefanucci, M. (2019). "Classification of Functional Fragments by Regularized Linear Classifiers with Domain Selection". In: *Biometrika* 106.1, pp. 161–180. ISSN: 14643510. DOI: 10.1093/biomet/asy060.
- Kurtek, S. and Drira, H. (2015). "A Comprehensive Statistical Framework for Elastic Shape Analysis of 3D Faces". In: *Computers and Graphics (Pergamon)* 51, pp. 52–59. ISSN: 00978493. DOI: 10.1016/j.cag.2015.05.027.
- Lee, S., Charon, N., Charlier, B., Popuri, K., Lebed, E., Sarunic, M. V., Trouvé, A., and Beg, M. F. (2017). "Atlas-Based Shape Analysis and Classification of Retinal Optical Coherence Tomography Images Using the Functional Shape (Fshape) Framework". In: *Medical Image Analysis* 35, pp. 570–581. ISSN: 13618423. DOI: 10.1016/j.media.2016.08.012.
- Li, J. J. and Tong, X. (2020). "Statistical Hypothesis Testing versus Machine Learning Binary Classification: Distinctions and Guidelines". In: *Patterns* 1.7, p. 100115. ISSN: 26663899. DOI: 10.1016/j.patter.2020.100115.
- Lila, E. and Aston, J. A. D. (2020). "Statistical Analysis of Functions on Surfaces, With an Application to Medical Imaging". In: *Journal of the American Statistical Association* 115.531, pp. 1420–1434. ISSN: 0162-1459, 1537-274X. DOI: 10.1080/01621459.2019.1635479.
- Lila, E., Aston, J. A., and Sangalli, L. M. (2016). "Smooth Principal Component Analysis over Two-Dimensional Manifolds with an Application to Neuroimaging". In: *Annals of Applied Statistics* 10.4, pp. 1854–1879. ISSN: 19417330. DOI: 10.1214/16-AOAS975.
- Lin, L., St. Thomas, B., Zhu, H., and Dunson, D. B. (2017). "Extrinsic Local Regression on Manifold-Valued Data". In: *Journal of the American Statistical Association* 112.519, pp. 1261–1273. ISSN: 1537274X. DOI: 10.1080/01621459.2016.1208615.

- Mai, Q., Zou, H., and Yuan, M. (2012). "A Direct Approach to Sparse Discriminant Analysis in Ultra-High Dimensions". In: *Biometrika* 99.1, pp. 29–42. ISSN: 00063444. DOI: 10.1093/biomet/asr066.
- Marron, J. and Dryden, I. L. (2021). *Object Oriented Data Analysis*. 1st ed. Boca Raton: Chapman and Hall/CRC. ISBN: 978-1-351-18967-5. DOI: 10.1201/9781351189675.
- Mejia, A. F., Yue, Y., Bolin, D., Lindgren, F., and Lindquist, M. A. (2020). "A Bayesian General Linear Modeling Approach to Cortical Surface fMRI Data Analysis". In: *Journal of the American Statistical Association* 115.530, pp. 501–520. ISSN: 1537274X. DOI: 10.1080/01621459.2019.1611582.
- Miller, M. I., Trouvé, A., and Younes, L. (2006). "Geodesic Shooting for Computational Anatomy". In: *Journal of Mathematical Imaging and Vision* 24.2, pp. 209–228. ISSN: 09249907. DOI: 10.1007/s10851-005-3624-0.
- Müller, H. G. (2005). "Functional Modelling and Classification of Longitudinal Data". In: Scandinavian Journal of Statistics 32.2, pp. 223–240. ISSN: 03036898. DOI: 10.1111/j.1467-9469.2005.00429. x.
- Paige, C. C. and Saunders, M. A. (1982). "Algorithm 583: LSQR: Sparse Linear Equations and Least Squares Problems". In: *ACM Transactions on Mathematical Software (TOMS)* 8.2, pp. 195–209. ISSN: 15577295. DOI: 10.1145/355993.356000.
- Park, J., Ahn, J., and Jeon, Y. (2021). "Sparse Functional Linear Discriminant Analysis". In: *Biometrika*. ISSN: 0006-3444. DOI: 10.1093/biomet/asaa107.
- Pinelis, I. (2007). "Optimum Bounds for the Distributions of Martingales in Banach Spaces". In: *The Annals of Probability* 22.4, pp. 347–370. ISSN: 0091-1798. DOI: 10.1214/aop/1176988477.
- Preda, C. (2007). "Regression Models for Functional Data by Reproducing Kernel Hilbert Spaces Methods". In: *Journal of Statistical Planning and Inference* 137.3, pp. 829–840. ISSN: 03783758. DOI: 10.1016/j.jspi.2006.06.011.
- Quarteroni, A. (2009). "Numerical Models for Differential Problems". In: *Modeling, Simulation and Applications* 2, pp. 1–593. ISSN: 20375263. DOI: 10.1007/978-88-470-1071-0.
- Ramsay, J. O. and Silverman, B. W. (2015). Functional Data Analysis. New York: Springer-Verlag. ISBN: 978-0-08-097087-5. DOI: 10.1016/B978-0-08-097086-8.42046-5.
- Reimherr, M., Sriperumbudur, B., and Taoufik, B. (2018). "Optimal Prediction for Additive Function-on-Function Regression". In: *Electronic Journal of Statistics* 12.2, pp. 4571–4601. ISSN: 19357524. DOI: 10.1214/18-EJS1505.
- Reuter, M., Wolter, F. E., and Peinecke, N. (2006). "Laplace-Beltrami Spectra as 'Shape-DNA' of Surfaces and Solids". In: *CAD Computer Aided Design* 38.4, pp. 342–366. ISSN: 00104485. DOI: 10.1016/j.cad.2005.10.011.
- Sabuncu, M. R. et al. (2011). "The Dynamics of Cortical and Hippocampal Atrophy in Alzheimer Disease". In: *Archives of Neurology* 68.8, pp. 1040–1048. ISSN: 00039942. DOI: 10.1001/archneurol. 2011.167.
- Shin, H. (2008). "An Extension of Fisher's Discriminant Analysis for Stochastic Processes". In: *Journal of Multivariate Analysis* 99.6, pp. 1191–1216. ISSN: 0047259X. DOI: 10.1016/j.jmva.2007.08.001.
- Smith, S. M. et al. (2013). "Functional Connectomics from Resting-State fMRI". In: Trends in Cognitive Sciences 17.12, pp. 666–682. ISSN: 13646613. DOI: 10.1016/j.tics.2013.09.016.
- Su, J., Kurtek, S., Klassen, E., and Srivastava, A. (2014). "Statistical Analysis of Trajectories on Riemannian Manifolds: Bird Migration, Hurricane Tracking and Video Surveillance". In: *Annals of Applied Statistics* 8.1, pp. 530–552. ISSN: 19417330. DOI: 10.1214/13-AOAS701.
- Sun, X., Du, P., Wang, X., and Ma, P. (2018). "Optimal Penalized Function-on-Function Regression Under a Reproducing Kernel Hilbert Space Framework". In: *Journal of the American Statistical Association* 113.524, pp. 1601–1611. ISSN: 1537274X. DOI: 10.1080/01621459.2017.1356320.
- Tibshirani, R. (1996). "Regression Shrinkage and Selection via the Lasso". In: *Journal of the Royal Statistical Society. Series B (Methodological)* 58.1, pp. 267–288. ISSN: 0035-9246.

- Tong, H. and Ng, M. (2018). "Analysis of Regularized Least Squares for Functional Linear Regression Model". In: *Journal of Complexity* 49, pp. 85–94. ISSN: 10902708. DOI: 10.1016/j.jco.2018.08.001.
- Vaillant, M., Miller, M. I., Younes, L., and Trouvé, A. (2004). "Statistics on Diffeomorphisms via Tangent Space Representations". In: *NeuroImage* 23.SUPPL. 1, S161–S169. ISSN: 10538119. DOI: 10.1016/j.neuroimage.2004.07.023.
- Vaillant, M. and Glaunès, J. (2005). "Surface Matching via Currents". In: Lecture Notes in Computer Science. Ed. by G. Christensen and M. Sonka. Vol. 3565. Berlin, Heidelberg: Springer, pp. 381–392. ISBN: 0302-9743. DOI: 10.1007/11505730\_32.
- Vaillant, M., Qiu, A., Glaunès, J., and Miller, M. I. (2007). "Diffeomorphic Metric Surface Mapping in Subregion of the Superior Temporal Gyrus". In: *NeuroImage* 34.3, pp. 1149–1159. ISSN: 10538119. DOI: 10.1016/j.neuroimage.2006.08.053.
- Wachinger, C., Golland, P., Kremen, W., Fischl, B., and Reuter, M. (2015). "BrainPrint: A Discriminative Characterization of Brain Morphology". In: *NeuroImage* 109, pp. 232–248. ISSN: 10538119. DOI: 10.1016/j.neuroimage.2015.01.032.
- Wahba, G. (1990). Spline Models for Observational Data. Society for Industrial and Applied Mathematics. DOI: 10.1137/1.9781611970128.
- Wang, G. and Wang, Y. (2017). "Towards a Holistic Cortical Thickness Descriptor: Heat Kernel-Based Grey Matter Morphology Signatures". In: *NeuroImage* 147, pp. 360-380. ISSN: 10538119. DOI: 10.1016/j.neuroimage.2016.12.014.
- Wang, X. and Zhu, H. (2017). "Generalized Scalar-on-Image Regression Models via Total Variation". In: *Journal of the American Statistical Association* 112.519, pp. 1156–1168. ISSN: 1537274X. DOI: 10.1080/01621459.2016.1194846.
- Yao, F. and Müller, H. G. (2010). "Functional Quadratic Regression". In: *Biometrika* 97.1, pp. 49–64. ISSN: 00063444. DOI: 10.1093/biomet/asp069.
- Yeo, T. B. T. et al. (2011). "The Organization of the Human Cerebral Cortex Estimated by Intrinsic Functional Connectivity". In: *Journal of Neurophysiology* 106.3, pp. 1125–1165. ISSN: 00223077. DOI: 10.1152/jn.00338.2011.
- Younes, L. (2019). Shapes and Diffeomorphisms. Vol. 171. Berlin, Heidelberg: Springer. ISBN: 978-3-642-12054-1. DOI: 10.1007/978-3-642-12055-8.
- Yu, S., Wang, G., Wang, L., and Yang, L. (2021). "Multivariate Spline Estimation and Inference for Image-on-Scalar Regression". In: *Statistica Sinica* 31.3, pp. 1463–1487. ISSN: 10170405. DOI: 10.5705/ss.202019.0188.
- Yuan, M. and Cai, T. T. (2010). "A Reproducing Kernel Hilbert Space Approach to Functional Linear Regression". In: *Annals of Statistics* 38.6, pp. 3412–3444. ISSN: 00905364. DOI: 10.1214/09-A0S772.
- Zaetz, J. and Kurtek, S. (2015). "A Novel Riemannian Framework for Shape Analysis of Annotated Surfaces". In: Proceedings of the Proceedings of the 1st International Workshop on DIFFerential Geometry in Computer Vision for Analysis of Shapes, Images and Trajectories 2015. British Machine Vision Association, pp. 3.1–3.11. ISBN: 1-901725-56-1. DOI: 10.5244/c.29.diffcv.3.
- Zhang, Z., Wu, Y., Xiong, D., Ibrahim, J. G., Srivastava, A., and Zhu, H. (2022). "LESA: Longitudinal Elastic Shape Analysis of Brain Subcortical Structures". In: *Journal of the American Statistical Association*, pp. 1–15. ISSN: 0162-1459, 1537-274X. DOI: 10.1080/01621459.2022.2102984.
- Zhu, H., Li, T., and Zhao, B. (2023). "Statistical Learning Methods for Neuroimaging Data Analysis with Applications". In: *Annual Review of Biomedical Data Science* 6.1, pp. 73–104. ISSN: 2574-3414, 2574-3414. DOI: 10.1146/annurev-biodatasci-020722-100353.