

Projected Statistical Methods for Distributional Data on the Real Line with the Wasserstein Metric

Matteo Pegoraro*, Mario Beraha‡

November 30, 2021

Abstract

We present a novel class of *projected* methods to perform statistical analysis on a data set of probability distributions on the real line, with the 2-Wasserstein metric. We focus in particular on Principal Component Analysis (PCA) and regression. To define these models, we exploit a representation of the Wasserstein space closely related to its weak Riemannian structure, by mapping the data to a suitable linear space and using a metric projection operator to constrain the results in the Wasserstein space. By carefully choosing the tangent point, we are able to derive fast empirical methods, exploiting a constrained B-spline approximation. As a byproduct of our approach, we are also able to derive faster routines for previous work on PCA for distributions. By means of simulation studies, we compare our approaches to previously proposed methods, showing that our *projected* PCA has similar performance for a fraction of the computational cost and that the *projected* regression is extremely flexible even under misspecification. Several theoretical properties of the models are investigated, and asymptotic consistency is proven. Two real world applications to Covid-19 mortality in the US and wind speed forecasting are discussed.

Keywords: Wasserstein spaces, Wasserstein metric, Principal Component Analysis, Wasserstein Regression, Metric Projection, Monotonic B-splines

1. Introduction

In many fields of machine learning and statistics, performing inference on a set of distributions is an ubiquitous but arduous task. The Wasserstein distance provides a powerful tool to compare distributions, as it requires very little assumptions on them and is at the same time reasonably easy to compute numerically. In fact, many other distances for distributions either require the existence of a probability density function or are impossible to evaluate, cf. Cuturi (2013), Peyré et al. (2019), Panaretos and Zemel (2020).

The Wasserstein distance recently gained popularity both in the statistics and machine learning community. See for instance Bassetti et al. (2006), Bernton et al. (2019a), Catalano et al. (2021) for statistical properties of the Wasserstein distance, Cao et al. (2019), Cuturi et al. (2019) and Cuturi and Doucet (2014) for applications in the field of machine and deep learning, Bernton et al. (2019b) and Srivastava et al. (2015) for applications in Bayesian computation.

In this work, we focus on the situation in which the single observation itself can be seen as a distribution, as in the analysis of images (Cuturi and Doucet, 2014; Banerjee et al., 2015), census data (Cazelles et al., 2018), econometric surveys Potter et al. (2017)

*. MOX – Department of Mathematics, Politecnico di Milano

†. Department of Mathematics, Politecnico di Milano

‡. Department of Computer Science, Università degli Studi di Bologna

and process monitoring (Hron et al., 2014). In particular, we consider observations to be distributions on the real line. There exist several possible ways to represent distributions, such as histograms, probability density functions (pdfs) and cumulative density functions (cdfs), each characterized by different constraints. For instance, histograms sum to one, pdfs integrate to one, and the limits for cdfs are 0 and 1, moreover all of these functions are nonnegative. These constraints translate into complex geometrical structures that characterize the underlying spaces these objects live in.

1.1 Previous work on distributional data analysis

One of the first works defining PCA for a data set of distributions is Kneip and Utikal (2001), where the authors apply tools from functional data analysis (FDA) directly to a collection of probability density functions. This approach, however, completely ignores the constrained nature of probability density functions, leading to poor interpretability of the results.

Based on theoretical results in Egozcue et al. (2006), who defines a Hilbert structure on a space of probability density functions on a compact interval (called a Bayes space), Delicado (2011) and Hron et al. (2014), propose a more reasonable approach to the problem of PCA for density functions. In particular, in Hron et al. (2014), the authors use the geometric properties of the Bayes space, coupled with a suitable transformation from the Bayes space to an L_2 space, to perform PCA on a set of pdfs using FDA tools, and then map back the results to the Bayes space.

Another, perhaps less widely used, approach focuses on borrowing tools from symbolic data analysis (SDA) in the context of histogram data (Nagabhushan and Pradeep Kumar, 2007; Rodríguez et al., 2000; Le-Rademacher and Billard, 2017). Moreover, in Verde et al. (2015) some of these attempts are extended to generic distributional data using Wasserstein metrics.

Finally, Bigot et al. (2017) and Cazelles et al. (2018) propose two PCA formulations based on the geometric structure of the Wasserstein space: a *geodesic* PCA and a *log* PCA. In a similar fashion, the recent preprints of Chen et al. (2021) and Zhang et al. (2020) propose linear regression and autoregressive models, respectively, for distributional data using the Wasserstein geometry.

We now highlight some key aspects of the aforementioned approaches. Hron et al. (2014) assumes that all the probability measures have the same support. This is hardly verified in practice, so that to apply their techniques one needs either to truncate the support of some of the probability density functions, or to extend others (for instance, by adding a small constant value and renormalizing), leading to numerical instability as discussed in Sections 7 and 8.

The SDA-based methods in Nagabhushan and Pradeep Kumar (2007); Rodríguez et al. (2000); Le-Rademacher and Billard (2017) and Verde et al. (2015) share the poor interpretability of SDA.

The methods in Bigot et al. (2017), Cazelles et al. (2018), Chen et al. (2021) and Zhang et al. (2020) are based on the weak Riemannian structure of the Wasserstein space, cf. Section 2.2. Such structure enables the authors to borrow ideas and terminologies from statistical frameworks defined on Riemannian manifolds (see Bhattacharya et al., 2012; Pennec, 2006, 2008; Huckemann et al., 2010; Patrangenaru and Ellingson, 2015; Fletcher, 2013; Banerjee et al., 2015). We can roughly distinguish those frameworks in two main approaches: the intrinsic/geodesic one and extrinsic/log one.

Briefly, intrinsic methods are defined using the metric structure of the Wasserstein space, working with geodesic curves and geodesic subsets, so that they faithfully respect

the metric of the underlying space. However, in general, intrinsic methods present many practical difficulties in that the optimization problems they lead to are usually nontrivial, as we discuss in Section 5.3. Instances of intrinsic methods for distributional data are the *geodesic* PCA in Bigot et al. (2017) and, under some rather restrictive assumptions, the linear models in Chen et al. (2021) and the autoregressive models in Zhang et al. (2020), see Sections 3.3 and 3.4.

On the other hand, extrinsic methods resort to the linear structure of suitably defined tangent spaces, by mapping data from the Wasserstein space to the tangent (through the so-called *log* map) and then mapping back the results to the Wasserstein space (through the *exp* map). Of course, this approach is less respectful of the underlying geometry than the intrinsic one, but usually presents several numerical advantages. An example of such extrinsic methods defined in the Wasserstein space is the *log* PCA in Cazelles et al. (2018).

The main issue with this *log* PCA is that the image of the *log* map inside the tangent of the Wasserstein space is not a linear space, but rather a convex cone embedded in a linear space (see Section 2.2). Hence, while exploiting the linear structure of the tangent, it is possible that the projection of some points onto the principal components end up outside of the cone. For these points, the *exp* map from the tangent to the Wasserstein space used in Cazelles et al. (2018) is not a metric projection, which in general is not available, so that the results in this setting are hardly interpretable.

1.2 Our contribution and outline

The contribution of this work is three folded. First, we propose alternative PCA and regression models for distributional data in the Wasserstein space. We term these models *projected*, in opposition to the *log* PCA in Cazelles et al. (2018). Second, by exploiting a geometric characterization of Wasserstein space closely related to its weak Riemannian structure, we build a novel approximation of the Wasserstein space using monotone B-spline. This allows us to represent the space of probability measures as a convex polytope in \mathbb{R}^J . Lastly, we obtain faster optimization routines for the *geodesic* PCAs defined in Bigot et al. (2017), exploiting the aforementioned B-spline representation.

Our *projected* framework lies in between the *log* one and the *geodesic* one, since we use an analogous to the *log* map to transform our data, as for extrinsic methods, but do not resort to the *exp* map to return to the Wasserstein space, using instead the metric projection operator. Thanks to this, our *projected* methods are more respectful of the underlying geometry than the *log* ones, while at the same time retaining the same reduced computational complexity. Thus, the *projected* methods expand the range of situations where *extrinsic* methods are an effective and efficient alternative to intrinsic tools: in our examples, the performance loss in general is marginal (see Section 7).

By centering the analysis in appropriate points of the Wasserstein space, one can identify the space of probability measures (with finite second moment) with the space of square integrable monotonically non-decreasing functions on a compact set. We use a suitable quadratic B-spline expansion to get a very handy representation of such functions. Through such B-spline expansion, it is possible to approximate the metric projection onto the Wasserstein space as a constrained quadratic optimization problem over a convex polytope, that is a well-established problem, cf. Potra and Wright (2000). This allows us to exploit the underlying linear structure of an L_2 space, so that all the machinery developed for functional data analysis can be directly applied to this setting. We address the issue of interpretability of the results, tackling a number of diverse applications and developing different ways to measure the loss of information caused by the *extrinsic* nature of our methods.

We observe that the idea of representing nondecreasing functions through B-splines for statistical purposes has been proposed also by Das and Ghosal (2017), in the context of Bayesian quantile regression, where the authors use B-splines with (random) monotonic coefficients as a generative model for random quantile functions. However, their focus is on defining a generative model, and not on developing a statistical setting exploiting the geometry given by the constrained representation. Along this direction, they do not restrict their attention to quadratic splines and consider cubic ones.

As already mentioned, a further contribution of this work is the derivation of alternative numerical optimization schemes for the *geodesic* PCA in Bigot et al. (2017) and Cazelles et al. (2018), based on the proposed quadratic B-spline expansion.

The remaining of the paper is organized as follows. Section 2 covers the basic concepts of Wasserstein distance and the weak Riemannian structure of the Wasserstein space, along with a brief discussion on a suitable way to exploit such structure for our purposes. Section 3 defines the *projected* PCA and *projected* regression in a general setting. In Section 4 we discuss the choice of the base point in which we center our analysis and how to efficiently approximate the metric projection through B-splines; in Section 5 we present the numerical algorithms needed to compute our *projected* methods and an alternative optimization routine for the *geodesic* PCA in Cazelles et al. (2018). Section 6 discusses the asymptotic properties of the spline approximation and of the *projected* models, establishing consistency of the estimators under some assumptions. Numerical illustrations on real and simulated data sets are shown in Sections 7 and 8. In particular, we apply our projected methods to two real world problems: we perform PCA on the US data on Covid-19 mortality by age and sex and perform a distribution regression to forecast the wind speed near a wind farm. Finally, the article concludes in Section 9. The Appendix collects all the proofs of the theoretical results, additional details on the simplicial PCA and regression, and further simulations. Code for reproducing the numerical results is available at <https://github.com/mberaha/ProjectedWasserstein>.

2. Preliminaries

In the following, we will consider probability measures on the real line \mathbb{R} endowed with the usual Borel σ -field, we will skip references to the σ -field whenever it is obvious.

Given a measure μ on \mathbb{R} define its cumulative distribution function $F_\mu(x) = \mu((-\infty, x])$ for $x \in \mathbb{R}$ and the associated quantile function $F_\mu^-(t) = \inf\{x \in \mathbb{R} : t \leq F_\mu(x)\}$. When F_μ is continuous and strictly monotonically increasing, $F_\mu^- = (F_\mu)^{-1}$.

2.1 Wasserstein metric and Wasserstein spaces

We start by recalling the definition of the 2-Wasserstein distance between two probability measures μ, ν on \mathbb{R} :

$$W_2^2(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \int_{\mathbb{R} \times \mathbb{R}} |x - y|^2 d\gamma(x, y), \quad (1)$$

where $\Gamma(\mu, \nu)$ is the collection of all probability measures on $\mathbb{R} \times \mathbb{R}$ with marginals μ and ν . Closely related to the definition of Wasserstein distance lies the one of Optimal Transport (OT). In particular, (1) identifies the Wasserstein distance with the minimal total transportation cost between μ and ν in the Kantorovich problem with quadratic cost (Ambrosio et al., 2008).

For our purposes, it is convenient to consider another formulation of the OT problem, originally introduced in Monge (1781). Given two measures μ, ν as before, the optimal

transport map from μ to ν is the solution of the problem

$$\inf_{T:T\#\mu=\nu} \int_{\Omega} |x - T(x)|^2 d\mu(x), \quad (2)$$

where $\#$ denotes the pushforward operator, that is for any measurable set B and measurable function

$$f : \mathbb{R} \rightarrow \mathbb{R}, \quad (f\#\mu)(B) = \mu(f^{-1}(B)). \quad (3)$$

Note that any solution of (2) induces one and only one solution of (1); moreover if the OT problem has a unique solution, then also the Wasserstein distance problem has only one solution. However not all Wasserstein distance problems can be solved through Monge's formulation (Ambrosio et al., 2008).

The unidimensional setting is a remarkable exception in that there exist explicit formulas for both problems. In particular, the Wasserstein distance can be computed as

$$W_2^2(\mu, \nu) = \int_0^1 |F_{\mu}^{-}(s) - F_{\nu}^{-}(s)|^2 ds, \quad (4)$$

and, if the measure μ has no atoms, then there exists a unique solution to Monge's problem given by $T_{\mu}^{\nu} = F_{\nu}^{-} \circ F_{\mu}$. For a proof of these results, see Chapter 6 of Ambrosio et al. (2008).

It is clear that, in general, the Wasserstein distance between two probability measures can be unbounded (for instance when in (4) F_{μ}^{-} is not square integrable on $[0, 1]$). Nonetheless, when restricting the focus on the set of probability measures with finite second moment, then it holds that W_2 defines a metric (see, for instance, Chapter 7 of Villani, 2008). Formally, let the Wasserstein space:

$$\mathcal{W}_2(\mathbb{R}) = \left\{ \mu \in \mathcal{P}(\mathbb{R}) : \int_{\mathbb{R}} x^2 d\mu < +\infty \right\}$$

then $(\mathcal{W}_2(\mathbb{R}), W_2)$ is a separable complete metric space.

2.2 Weak Riemannian structure of the Wasserstein Space

Thanks to the uniqueness of the transport maps, by fixing an absolutely continuous (a.c.) probability measure $\mu \in \mathcal{W}_2(\mathbb{R})$, we can associate to any $\nu \in \mathcal{W}_2(\mathbb{R})$ the optimal transport map T_{μ}^{ν} . Since $\int_{\mathbb{R}} |T_{\mu}^{\nu}(x)|^2 d\mu = \int_{\mathbb{R}} x^2 d\nu$ we can define the following map $\varphi_{\mu} : \mathcal{W}_2(\mathbb{R}) \rightarrow L_2^{\mu}(\mathbb{R})$ with the rule: $\varphi_{\mu}(\nu) = T_{\mu}^{\nu}$.

We note several immediate but interesting properties of the map φ_{μ} . First, it is an isometry (and so a homeomorphism onto its image) since

$$\int_{\mathbb{R}} |T_{\mu}^{\nu}(x) - T_{\mu}^{\eta}(x)|^2 d\mu = \int_{[0,1]} |F_{\nu}^{-} - F_{\eta}^{-}|^2 ds = W_2^2(\nu, \eta).$$

Second, the image of φ_{μ} is a closed convex cone in $L_2^{\mu}(\mathbb{R})$: a set closed under addition and positive scalar multiplication. In fact, for any $\lambda \geq 0$, λT_{μ}^{ν} is still a transport map from μ to another measure whose quantile is λF_{ν}^{-} ; and similarly $T_{\mu}^{\nu} + T_{\mu}^{\eta} = (F_{\nu}^{-} + F_{\eta}^{-}) \circ F_{\mu}$. Being $\mathcal{W}_2(\mathbb{R})$ complete, $\varphi_{\mu}(\mathcal{W}_2(\mathbb{R}))$ is closed in $L_2^{\mu}(\mathbb{R})$. Third, $\varphi_{\mu}(\mu) = id_{\mathbb{R}}$ (where id_C denotes the identity map of the set C). Finally, as shown in Panaretos and Zemel (2020), φ_{μ} is not surjective and $\varphi_{\mu}(\mathcal{W}_2(\mathbb{R}))$ is the set of μ -a.e. non decreasing functions in $L_2^{\mu}(\mathbb{R})$.

The inverse of the map of φ_{μ} is the measure pushforward (see Equation 3) and it is defined on the whole $L_2^{\mu}(\mathbb{R})$: given $f \in L_2^{\mu}(\mathbb{R})$, then $\nu = f\#\mu$ is a measure in $\mathcal{W}_2(\mathbb{R})$. In fact:

$$\int |x|^2 d\nu = \int |f(x)|^2 d\mu = \|f\|_{\mu}^2$$

A natural way to define a tangent structure for $\mathcal{W}_2(\mathbb{R})$ is therefore to take advantage of the cone structure given by φ_μ . In fact for closed convex cones, there are already notions of tangent cones. Similarly to Rockafellar and Wets (1998), Theorem 6.9, we can define:

$$\text{Tan}_\mu(\mathcal{W}_2(\mathbb{R})) := \text{Tan}_{id_{\mathbb{R}}}(L_2^\mu(\mathbb{R})) = \overline{\{f \in L_2^\mu(\mathbb{R}) \mid \exists h > 0 : id + hf \in \varphi_\mu(\mathcal{W}_2(\mathbb{R}))\}}^{L_2^\mu(\mathbb{R})} \quad (5)$$

We remark that Theorem 6.9 in Rockafellar and Wets (1998) is stated in \mathbb{R}^n , but it holds also more generally, for instance in an Hilbert space (see Aubin and Frankowska (2009), Chapter 4).

A geometric interpretation of (5) is the following. The tangent space consists of all the vectors f that move the base point inside the cone $\varphi_\mu(\mathcal{W}_2(\mathbb{R}))$, when considered up to a scale factor h . Hence, f plays the role of direction of a tangent vector going out from the tangent point. Furthermore, since for every $f \in \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$ then $f + id \in \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$ we have that $\varphi_\mu(\mathcal{W}_2(\mathbb{R}))$ is included in the tangent space. As shown later in this Section, the inclusion is strict and the tangent space is much larger than $\varphi_\mu(\mathcal{W}_2(\mathbb{R}))$.

Note that we can recover the definition of tangent space given by Ambrosio et al. (2008) and Panaretos and Zemel (2020) by a simple “change of variable”: calling $g = id + hf$ then substituting $(g - id)/h$ in (5) gives the following definition of tangent

$$\text{Tan}_\mu(\mathcal{W}_2(\mathbb{R})) = \overline{\{\lambda(f - id) \mid f \in \varphi_\mu(\mathcal{W}_2(\mathbb{R})); \lambda > 0\}}^{L_2^\mu(\mathbb{R})},$$

which is the one given in Ambrosio et al. (2008) and Panaretos and Zemel (2020). As shown in Panaretos and Zemel (2020) the tangent cone $\text{Tan}_\mu(\mathcal{W}_2(\mathbb{R}))$ is indeed a linear space. For this reason we refer to it as tangent space, instead of cone.

In analogy to Riemannian geometry, following Ambrosio et al. (2008) and Panaretos and Zemel (2020), we define the \log_μ and \exp_μ maps. Having fixed μ absolutely continuous:

$$\begin{aligned} \log_\mu : \mathcal{W}_2(\mathbb{R}) &\rightarrow \text{Tan}_\mu(\mathcal{W}_2(\mathbb{R})) & \exp_\mu : \text{Tan}_\mu(\mathcal{W}_2(\mathbb{R})) &\rightarrow \mathcal{W}_2(\mathbb{R}) \\ \nu &\mapsto T_\mu^\nu - id & f &\mapsto (id + f)\#\mu \end{aligned} \quad (6)$$

We briefly highlight some properties of these maps; properties which immediately follows from the discussion above.

Remark 1 *The map \log_μ is defined on the whole space $\mathcal{W}_2(\mathbb{R})$. Moreover, it is clearly an isometry: $W_2(\eta, \nu) = \|\log_\mu(\eta) - \log_\mu(\nu)\|_{L_2^\mu(\mathbb{R})}$ (Panaretos and Zemel, 2020). This shows that there is no local-approximation issue when working in the tangent space, in contrast with the usual Riemannian manifold setting. There, the tangent space usually provides good approximation only in a neighborhood of the tangent point.*

Remark 2 *The map \log_μ is not surjective on Tan_μ , indeed its image $\text{Im}(\log_\mu)$ is a closed convex subset of $L_2^\mu(\mathbb{R})$ given by all the maps f such that $f + id \in \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$, that is, $f + id$ is μ -a.e. increasing. The restriction of \exp_μ on $\text{Im}(\log_\mu)$, henceforth denoted by $\exp_\mu|_{\log_\mu(\mathcal{W}_2(\mathbb{R}))}$, is an isometric homeomorphism and its inverse is \log_μ . In particular, we observe that $\log_\mu \circ \exp_\mu$ is not a metric projection in L_2^μ . That is, in general $\log_\mu \circ \exp_\mu(f) \neq \arg \min_{g \in \text{Im}(\log_\mu)} \|f - g\|_{L_2^\mu}$.*

2.3 Intrinsic and extrinsic methods in the Wasserstein space

As mentioned in Section 1.1, borrowing ideas from Riemannian geometry leads to discerning statistical methods on the Wasserstein space in the classes of *intrinsic* and *extrinsic* methods.

The Weak Riemannian structure presented in Section 2.2 provides a suitable environment for developing intrinsic methods. In fact, the geodesic structure of $\mathcal{W}_2(\mathbb{R})$ can be recovered through the linear structure of any $L_2^\mu(\mathbb{R})$ space through the isometry φ_μ . Pointwise interpolation of the transport maps coincide with the geodesic between measures. In other words, given μ a.c., the geodesic between ν and η is given by:

$$\gamma(t) = ((1-t) \cdot T_\mu^\nu + t \cdot T_\mu^\eta) \# \mu \quad (7)$$

Thus, such geodesic structure can be recovered in many different (but equivalent) ways, depending on μ .

On the other hand, Remark 1 motivates the development of extrinsic tools, since working in the image of \log_μ inside the tangent space Tan_μ is exactly like working in $\mathcal{W}_2(\mathbb{R})$. This is not common in Riemannian manifold framework, since usually the tangent space provides a good approximation only near to the tangent point. As a consequence, if in the general Riemannian manifold framework the choice of the tangent point μ is crucial (since results for extrinsic methods might be significantly altered for different choices of μ) when working with $\mathcal{W}_2(\mathbb{R})$ this is not the case.

To further motivate this key point, consider μ and ν a.c. measures; the maps $\log_\nu \circ (\exp_{\mu|\log_\mu(\mathcal{W}_2(\mathbb{R}))})$ and $\varphi_\nu \circ \varphi_\mu^{-1}$ are isometric homeomorphisms (as composition of isometries and homeomorphisms). In other words, they preserve distances and send border elements of $\log_\mu(\mathcal{W}_2(\mathbb{R}))$ or $\varphi_\mu(\mathcal{W}_2(\mathbb{R}))$ into border elements of $\log_\nu(\mathcal{W}_2(\mathbb{R}))$ and $\varphi_\nu(\mathcal{W}_2(\mathbb{R}))$, respectively, and the same with internal points (and so in particular, they preserve distances from any point to the border). In Chen et al. (2021), Bigot et al. (2017) and Zhang et al. (2020) μ is chosen as the barycentric measure \bar{x} of the observations $x_i \in \mathcal{W}_2(\mathbb{R})$. The discussion above implies that considering the tangent space at the Wasserstein barycenter \bar{x} and working on $\log_{\bar{x}}(x_i) = \log_{\bar{x}}(x_i) - \log_{\bar{x}}(\bar{x})$ is exactly the same as considering the tangent space at any μ a.c. and working on $\log_\mu(x_i) - \log_\mu(\bar{x})$ for our statistical purposes. So the choice of the tangent space from the theoretical point of view is completely arbitrary.

Moreover, centering the analysis in the barycenter presents a drawback when studying asymptotic properties of the models under consideration, since \bar{x} changes as the sample size grows. In Section 4.1 we propose to fix μ as the uniform measure on $[0, 1]$. This choice not only allows us to derive empirical methods that are extremely simple to implement, cf. Section 5, but also allows us to study asymptotic properties of the models in Section 6.2 without resorting to parallel transport, as done for instance in Chen et al. (2021).

2.4 Tangent vs. L_2^μ

Lastly, we briefly discuss the major differences between using a tangent space representation of $\mathcal{W}_2(\mathbb{R})$ and using the representation given by some φ_μ .

We recall that, for a fixed μ a.c., the two representations are indeed quite similar $\varphi_\mu(\nu) = T_\mu^\nu$, $\log_\mu(\nu) = T_\mu^\nu - id$; a priori one may prefer the tangent representation, because it already expresses data as vectors coming out of a point. Therefore, for instance, it might result practically more convenient to center the analysis in the barycenter and work on vectors, taking away any “data centering” issues. At the same time, also notational coherence with already existing methods might benefit from this choice.

However, especially when dealing with extrinsic techniques, we found slightly more practical to use the φ_μ representation in that it is more straightforward to represent $\varphi_\mu(\mathcal{W}_2(\mathbb{R}))$ compared to $\log_\mu(\mathcal{W}_2(\mathbb{R}))$: the first one can in fact be represented directly as the cone of the μ -a.e non-decreasing functions.

3. Projected Models in the Wasserstein Space

In this section, exploiting the embeddings given by φ_μ , we define a class of *projected* statistical methods to perform extrinsic analysis for data in the Wasserstein space.

To give a general framework, we do not restrict our attention to a particular φ_μ yet, even though in Section 4 we argue that a natural choice which allows an easier implementation of the empirical methods is letting μ be the uniform distribution on $[0, 1]$. Hence, for the sake of notation, we consider a generic case of data lying in a closed convex cone X inside a separable Hilbert space H . In our setting, H would be $L_2^\mu(\mathbb{R})$ and $X = \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$, for some $\mu \in \mathcal{W}_2(\mathbb{R})$ absolutely continuous.

3.1 Principal component analysis

We start by defining one of the main contributions of our work: the *projected* PCA. We recall that for an H -valued random variable \mathcal{X} , PCA is a well established technique and amounts to finding the eigenfunctions of the Karhunen-Loève expansion of the covariance operator of \mathcal{X} , see Ramsay (2004). Observe that any X -valued random variable can be considered as an H -valued one (by the inclusion map), so that a notion of PCA is already available.

When defining principal components, a key notion is the one of dimension of the principal component (PC). In this work, principal components will be closed convex subsets of H , and we will always define the dimension of a subset of H as the dimension of the smallest affine subset of H containing it. For a generic closed convex set $C \subset H$, let Π_C denote the metric projection onto C : $\Pi_C(x) := \arg \min_{c \in C} \|x - c\|$ and, for a set of vectors U , denote with $Sp(U)$ its linear span.

In what follows, we denote by x_0 the “center” of the PCA. For us, $x_0 = \mathbb{E}[\mathcal{X}]$, or its empirical counterpart. To have a well defined PCA, we always assume that x_0 belongs to the relative interior of the convex hull of the support of \mathcal{X} , see Appendix A for the definition of relative interior and further details. This is a rather technical hypothesis but it is not a restrictive one. For instance, it is always verified for empirical measures and when $X \subseteq \mathbb{R}^d$ and hence for our empirical methods, cf. Section 5.1.

Definition 1 (*Projected PCA*). *Given \mathcal{X} a random variable with values in $X \subset H$, let $U_k = \{w_1, \dots, w_k\}$ be its first k H -principal components centered in $x_0 = \mathbb{E}[\mathcal{X}]$. A (k, x_0) -projected principal component of \mathcal{X} is the biggest closed convex subset $U_X^{x_0, k}$ of X such that: (i) $x_0 \in U_X^{x_0, k}$, (ii) $\dim(U_X^{x_0, k}) = k$, and (iii) $U_X^{x_0, k} \subseteq \Pi_X(Sp(U_k))$.*

In other words, the projected principal component is obtained by approximating the span of the principal components found in H , with convex subsets in X . Note that the principal components in H might “capture” some variability which is not present when measuring distances inside X . In fact the projection of a point belonging to X onto a direction w_j might end up being outside X , see Section 3.3. However, as we will show in Section 7, in our examples the projected PCA behaves well and this issue does not seem to affect significantly the performance.

Remark 3 *Convex sets are essential in our analysis since, thanks to (7) convex sets in X are precisely the subsets of $\mathcal{W}_2(\mathbb{R})$ which are geodesically complete: the geodesic connecting any pair of points in the subset, is contained in the subset. Geodesic subsets are a natural generalization of linear spaces.*

Remark 4 *The metric projection of a linear subspace onto a convex subset can end up being a nonconvex set. In addition to that, while loosing convexity, the dimension of the metric projection of a convex subset can be bigger of the dimension of the original subset. A simple example where both cases happen is the projection of $y = -x$ onto $x, y \geq 0$ in \mathbb{R}^2 .*

We observe that inside a projected principal component, we have a preferential orthonormal basis given by the principal components in H ; for this reason we call $U_k = \{w_1, \dots, w_k\}$ *principal directions*.

Although it might seem impractical to find the projected component, the following Lemma provides a more convenient alternative characterization.

Lemma 1 *Let x_0 and $U_X^{x_0, k}$ be as in Definition 1, then $U_X^{x_0, k} = (x_0 + Sp(U_k)) \cap X$.*

Natural alternatives to Definition 1 would be, for instance, to let the projected principal directions (component) be the metric projection of w_1, \dots, w_k (the linear span of $\{w_1, \dots, w_k\}$) onto X , respectively. In the former case, the projection would not guarantee the orthogonality of the projected directions, which is instead essential to properly explore the variability. Moreover, since the “tip” of the projected unit vectors would likely lie on the border of X , the projection of a new observations on a direction would still lie outside of X as soon as the score associated to that direction is larger than 1. The latter case, instead, presents the drawbacks pointed out in Remark 4.

We argue that, despite its simplicity, Definition 1 is indeed very well suited for statistical analysis in the Wasserstein Space. For instance, we are guaranteed that, as the dimension grows up, the k projected components provide a monotonically better fit to the data. This is easily verified because Π_X is a strictly non-expansive operator, being X closed and convex (see Deutsch (2012)), which implies the following Proposition.

Proposition 1 *With the same notation as Definition 1, for any $x \in X$ we have:*

$$\|\Pi_{U_X^{x_0, k}}(x) - x\| \geq \|\Pi_{U_X^{x_0, k+1}}(x) - x\| \rightarrow 0 \text{ with } k \rightarrow +\infty.$$

Once a principal component is found, a classical task that one may want to perform is to project a new “observation” $x^* \in X$ onto $U_X^{x_0, k}$, for instance for dimensionality reduction purposes. In general, the metric projection on generic convex subsets might be arduous to find, we will deal with this issue in Section 4. Nevertheless, we can use the following Proposition to reduce in advance the dimension of the parameters involved in the problem; turning it into a projection problem inside the principal projected component, which allows for faster computations (see Equation 13).

Proposition 2 *Let $x^* \in X$ and let Π_k be the orthogonal projection on $Span(U_k)$. The projection of x^* onto $U_X^{x_0, k}$ is given by*

$$\arg \min_{v' \in U_X^{x_0, k}} \|x^* - v'\| = \Pi_{Sp(U_k) \cap (X - x_0)}(\Pi_k(x^* - x_0)) + x_0. \quad (8)$$

Lastly, we observe that, since projected principal components are not linear subspaces, the scores of some points on a principal direction can vary as we increase the dimension of the principal component.

3.2 Regression

Broadly speaking, a regression model between two variables with values in two different spaces is given by an operator between such spaces, which for every input value of the independent variable, returns a predicted value for the dependent variable. In the following, let us denote with \mathcal{Z} the independent variable and with \mathcal{Y} the dependent one. A regression model is usually understood as an operator Γ specifying the conditional value of \mathcal{Y} given \mathcal{Z} , that is, $\mathbb{E}[\mathcal{Y}|\mathcal{Z}] = \Gamma(\mathcal{Z})$.

If the spaces where \mathcal{Z} and \mathcal{Y} take values possess a linear structure, this linearity is usually exploited by means of a (kernel) linear operator, with possibly an “intercept” term. To define our *projected* regression model, we want to exploit the cone structure of X in a similar fashion. In fact, such linear kernel operators combine good optimization properties and interpretability since their kernels can provide insights into the analysis, much like coefficients in multivariate linear regression.

We treat separately the cases where the X -valued variable is the independent or the dependent one. The case when both variables are X -valued follows naturally. To keep the notation light, in what follows we will not distinguish between “proper” linear operators and linear operators with an added intercept term, which could as well be employed in all the incoming definitions to gain flexibility.

Consider the case in which we have an independent X -valued random variable, and denote with V the space where the dependent variable takes value. Despite the fact that X is not a linear space, with an abuse of notation, we call “linear” an operator which respect sum and positive scalar multiplication for elements in X . Such operators are in fact obtained by restricting on X linear operators defined on H . Following this idea, in order to define linear regression for an X -valued independent random variable, we consider such variable as H -valued, obtain the regression operator and then take the restriction of the operator on X . In this way, when $H = L_2^\mu(\mathbb{R})$ and $X = \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$, it is possible to exploit the classical FDA framework to perform all kinds of distribution on scalar/vector/etc... regression. For brevity, we report only the definition with $V = \mathbb{R}$.

Definition 2 *Let \mathcal{Z} an X -valued random variable, and \mathcal{Y} a real valued one. Let $\Gamma_\beta : H \rightarrow \mathbb{R}$ be a functional linear regression model for such variables, with \mathcal{Z} considered as H -valued and $\Gamma_\beta(v) = \langle \beta, v \rangle$. A projected linear regression model for $(\mathcal{Z}, \mathcal{Y})$ is given by $(\Gamma_\beta)|_X$.*

Now we turn to the cases which feature an X valued independent variable and a Z valued dependent one, for Z a generic Hilbert space. Through the inclusion $X \hookrightarrow H$, we can consider a regression problem with X -valued dependent variable, as a problem with H -valued dependent variable. Comparing this situation with the previous one, it is clear that we now face a “dual” problem. Indeed, while before we needed to restrict the domain from H to X , we now need to force the codomain of Γ to lie inside X . We would like to retain the same properties that make linear kernel operators appealing as regression operators between Hilbert spaces. A possibility could be considering a linear kernel operator Γ with values in H and restricting it to $\Gamma^{-1}(X)$. However, this would imply that for any $z \notin \Gamma^{-1}(X)$ no prediction would be available.

We argue that a more reasonable approach consists in finding an operator $\Gamma_P : Z \rightarrow X$ as close as possible (in some sense that will be clear later) to the linear kernel operator Γ aforementioned. Hence, we relax the linearity assumption in favor of Lipschitzianity, and take as regression operator $\Pi_X \circ \Gamma$, whose image always lies in X . Note that Γ_P inherits the interpretability of the kernel of Γ .

To motivate such choice, we give the following notion of a projected operator.

Definition 3 Let Z be a normed space and consider \mathcal{Z} a Z -valued random variable. Let $\Gamma : Z \rightarrow H$ a generic Lipschitz operator between Z and H . A (\mathcal{Z}, X) -projection of Γ is an operator $\Gamma_{\mathcal{P}} : Z \rightarrow X$ such that:

$$\Gamma_{\mathcal{P}} = \arg \min_{T:Z \rightarrow X} \mathbb{E}_{\mathcal{Z}}[\|\Gamma(v) - T(v)\|^2]$$

In other words, $\Gamma_{\mathcal{P}}$ provides the best pointwise approximation of the H -valued operator Γ , averaged w.r.t. the measure induced by \mathcal{Z} . Hence, given a \mathcal{Z} a Z -valued random variable and \mathcal{Y} an X -valued random variable and a linear regression model $\Gamma : Z \rightarrow H$ for $(\mathcal{Z}, \mathcal{Y})$, the projected regression model induced by Γ is $\Gamma_{\mathcal{P}}$.

Proposition 3 With the same notation as above, if $\mathbb{E}[\|\mathcal{Z}\|^2] < \infty$, then $\Gamma_{\mathcal{P}} = \Pi_X \circ \Gamma$.

Proof For any $T : Z \rightarrow X$, it holds: $\|\Gamma(z) - \Pi_X(\Gamma(z))\| \leq \|\Gamma(v) - T(v)\|$. Moreover, Γ and $\Pi_X \circ \Gamma$ are Lipschitz, and being Π_X non-expansive, they share the same constant $L > 0$:

$$\|\Gamma(v) - \Pi_X \circ \Gamma(v)\|^2 \leq 2L\|v\|^2$$

and thus $\mathbb{E}_{\mathcal{Z}}[\|\Gamma(z) - \Pi_X \circ \Gamma(z)\|^2]$ is bounded iff \mathcal{Z} has finite second moment. \blacksquare

The only case left out from the treatment above is when both the independent and the dependent variables are X -valued. This case, however, follows naturally by combining the two approaches and we report the definition below.

Definition 4 Let \mathcal{Z} and \mathcal{Y} two X -valued random variables. Let $\Gamma : H \rightarrow H$ be a functional linear regression model for the variables considered as H -valued. A projected linear regression model for $(\mathcal{Z}, \mathcal{Y})$ is given by $(\Pi_X \circ \Gamma)|_X$.

Remark 5 When considering a regression with X -valued independent variable, one may want to relax the restriction on X in Definition 2 for various reasons; for instance one may have measurement errors, or by design the test set may consider points also outside X . In such cases it is worth considering the problem of how many continuous linear extensions of $\Gamma|_X$ are possible on the whole H . A sufficient condition for the uniqueness of such extension is the following: there exist a sequence of linear subspaces of H , say $\{H_J\}_{J \geq 1}$, such that $\bigcup_J H_J$ is dense in H and $X_J := H_J \cap X$ contains a basis of H_J for every J .

Remark 6 When $H = L_2^\mu(\mathbb{R})$ and $X = \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$ the condition in Remark 5 is verified, for instance, by Remark 8 in Section 4.3. Moreover, observe that the uniqueness of the extension can also be proven thank to Jordan's representation of functions $f : \mathbb{R} \rightarrow \mathbb{R}$ with bounded variation (BV). In fact any f with BV can be written as the difference of monotone functions and thus $\Gamma(f)$ is fixed. Then by the density of BV functions in H , we define Γ on the remaining elements of H .

3.3 Comparison with intrinsic methods

We now compare the projected methods defined earlier in this Section and the intrinsic counterparts. In particular, we focus on the *geodesic* PCA defined in Bigot et al. (2017) and Cazelles et al. (2018) and on the distribution on distribution regression model in Chen et al. (2021).

Bigot et al. (2017) and Cazelles et al. (2018) define two different PCA, namely a global and a nested one; in particular the nested approach presents analogies with other PCAs developed for manifold valued random variables (Jung et al., 2012; Huckemann and Eltzner, 2018; Pennec, 2018); we report the two definitions below.

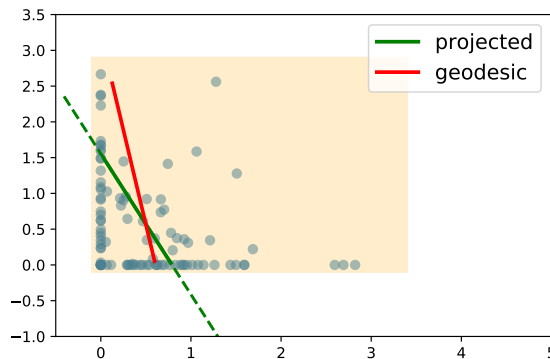


Figure 1: Comparison of projected and geodesic PCA when $H = \mathbb{R}^2$ and X is the shaded rectangle. The projected principal direction is rather different from the geodesic one because most of the observations (blue dots) are concentrated around the borders

Definition 5 (*Global geodesic PCA*) Let \mathcal{X} a random variable with values in X with $\mathbb{E}[\mathcal{X}] = x_0$. A (k, x_0) -global geodesic PC is a set C^* minimizing $\mathbb{E}[d(\mathcal{X}, C)^2]$ over the closed convex sets $C \subset X$ such that $x_0 \in C$ and $\dim(C) \leq k$

Definition 6 (*Nested geodesic PCA*) Let \mathcal{X} a random variable with values in X with $\mathbb{E}[\mathcal{X}] = x_0$. For $k = 1$, a (k, x_0) -nested geodesic PC is a set C_k^* such that C_k^* is a minimizer of $\mathbb{E}[d(\mathcal{X}, C)^2]$ over the closed convex sets $C \subset X$ such that $x_0 \in C$ and $\dim(C) \leq k$; for $k \geq 1$, a (k, x_0) -nested geodesic PC is a set C_k^* such that C_k^* is a minimizer of $\mathbb{E}[d(\mathcal{X}, C)^2]$ over the closed convex sets $C \subset X$ such that: $x_0 \in C$, $\dim(C) \leq k$, and $C \supset C_{k-1}^*$, where C_{k-1}^* is a $(k-1, x_0)$ -nested geodesic PC.

The first key difference between the global and the nested geodesic PCA is that the latter provides a notion of preferential directions in the principal component, while the first one does not. In fact, the first nested principal component corresponds to the first principal direction, and it is possible to find the remaining principal directions by imposing orthogonality constraints as we obtain nested PCs of higher dimensions. Thus, the nested geodesic PCA is more suitable to explore and visualize the variability in a data set, see also Section 7. On the other hand, exactly because of the lack of such constraints, the global PCA is in general more flexible and provides superior performance in terms of *reconstruction error*, cf. Section 7.

Comparing these definitions with the one of our projected PCA, the key difference is that geodesic PCAs do not exploit the Hilbert structure of H . Thus, as we discuss in Section 5.3, the numerical routines needed to find such principal components rely on nonlinear constrained optimization, which can be extremely demanding and nontrivial to implement. This is in sharp contrast with our projected PCA in Definition 1, that, thanks to Lemma 1 can be straightforwardly computed. However, as a result, the projected PCA is in general less respectful of the underlying metric structure. By investigating this issue in simpler settings, for instance when $H = \mathbb{R}^d$ and X is a convex polytope in \mathbb{R}^d , we noticed that the differences between the projected principal directions and the nested geodesic ones become appreciable only if the random variable \mathcal{X} gives significant probability to values near the borders of X . See for instance Figure 1. While this intuition remains valid also in the more complex setting that we investigate in this paper, it is harder to imagine realizations of \mathcal{X} near the borders of X .

Note that the interpretability of the projected PCA is determined by the level of discrepancy between the definitions, as in Figure 1, which depends on how much variability

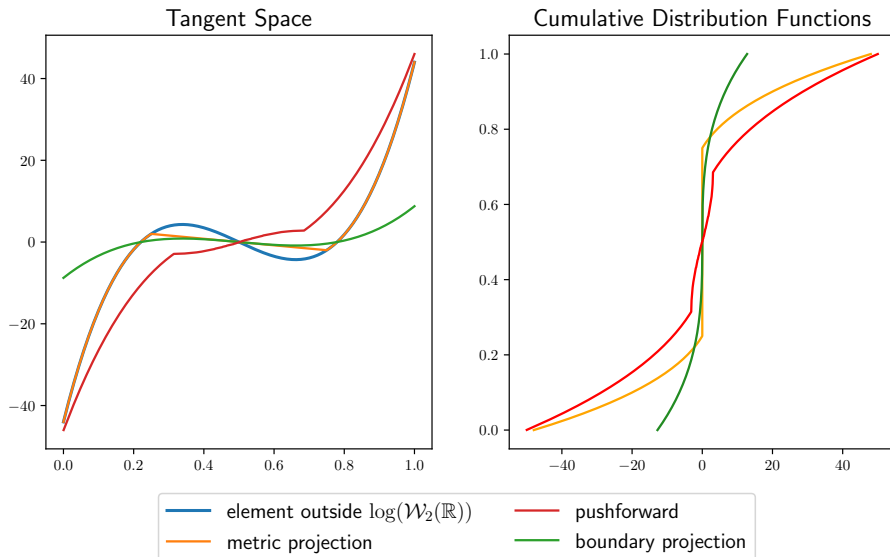


Figure 2: Comparison between different projections onto X for a point $x \in H \setminus X$ (blue line) in the tangent space (left panel) and the associated cumulative distribution functions (right panel) when the base point μ is the uniform measure on $[0, 1]$. The orange, green and red curves are obtained with metric projection, boundary projection and $\log_\mu \circ \exp_\mu$ respectively.

it is correctly captured by the component, that is how much of the variability captured by the projected component lies in X . This intuition is formalized in Section 7.2 where two measures of “reliability” of the projected PCA are proposed.

Turning to the regression context, Chen et al. (2021) define a distribution on distribution linear regression model in the Wasserstein space. Their approach considers two different tangent spaces of $\mathcal{W}_2(\mathbb{R})$ (the first one centered in the barycenter of the independent variable and the second one centered in the barycenter of the dependent variable) and map the observations to the corresponding tangent spaces. They then use FDA tools to estimate a functional linear model $\hat{\Gamma}$ between those two spaces. When the image of the regression operator Γ lies inside the image of the log map centered in the dependent variable’s barycenter, their distribution on distribution regression can be considered a properly intrinsic method. This assumption is used to prove asymptotic properties of their methodology, but as the authors in Chen et al. (2021) notice, is hardly verified in practice, so that whenever the output of the regression operator is not a distribution, they resort to squeezing such a value with some scalar multiplication, namely “boundary projection”, which in general is not a metric projection. The boundary projection step gives an extrinsic nature to their model and we provide further comparisons with our methods in Section 3.4.

3.4 Comparison with other extrinsic methods

In this section, we offer a comparison of our projected methods with other extrinsic methods, namely the log PCA in Cazelles et al. (2018) and the distribution on distribution regression in Chen et al. (2021), which, as outlined in the previous section, may behave as an extrinsic method. Let us start with the former.

Cazelles et al. (2018) propose the definition of a log PCA as an alternative to the geodesic PCAs in Bigot et al. (2017). Both the log and the projected PCA are extrinsic

methods: they proceed by carrying out the PCA in a linear space H and then map back the results to the Wasserstein space, following an approach which had already been proposed by Fletcher et al. (2004).

For the log PCA, H is the tangent space at μ , for the projected H is $L_2^\mu(\mathbb{R})$. Given $U_k = \{w_1, \dots, w_k\}$ the first k H -principal components, the log principal component in $\mathcal{W}_2(\mathbb{R})$ is $\exp_\mu(Sp(U_k))$. Analogously, by considering the convex cone $X =: \log_\mu(\mathcal{W}_2(\mathbb{R})) \subseteq H$, the principal component in X is $\log_\mu(\exp_\mu(Sp(U_k)))$.

We notice two key differences between the log and projected PCA. First, as pointed out in Remark 2, $\log_\mu \circ \exp_\mu$ is not a metric projection in L_2^μ so that given a point $x \in H \setminus X$, $\log_\mu(\exp_\mu(x))$ might end up being extremely different from x . See for instance Figure 2 where for a point x (blue line) that is close (in the L_2^μ norm) to X , $\log_\mu(\exp_\mu(x))$ turns out to be quite far from x . In the context of PCA, this means that as soon as the projection onto $Sp(U_k)$ of observation lies outside of X , the log PCA quickly loses its interpretability. Second, as discussed in Remark 4, there is no guarantee that $\log_\mu(\exp_\mu(Sp(U_k)))$ is contained in $Sp(U_k)$, its dimension might increase and it might not even be convex. For this same reason, in general, log PCA cannot define a set of (orthogonal) principal directions which span the principal component. Hence, it is not possible to work directly on the scores of the PCA.

Combined, we believe that the above mentioned issues present a major drawback of the log PCA when compared to the projected PCA, as they prevent the possibility of doing proper dimensionality reduction and working on the scores of data points on the principal components. Finally, we also point out that approximating the \exp_μ map is a nontrivial task, involving computing numerically the preimages of an arbitrary large number of sets and numerical differentiation, that can lead to numerical instability of the log PCA.

We end this discussion with a comparison between the boundary projection in Chen et al. (2021) and the metric projection. Their difference, for a possible regression output $x \in H \setminus X$ is depicted in Figure 2. Note that, by construction, such a procedure shrinks the tails of the output. Even when the regression output is slightly outside the image of the log map, the boundary projection result can be extremely far from the regression output and from the metric projection in terms of Wasserstein distance. For example, in Figure 2, the regression output and the projected method assign positive probability to values in the range $[-45, 45]$, while the output of the boundary projection assigns zero probability to values outside $[-17, 17]$. This underrepresentation of the variability might be a crucial issue depending on the application considered.

4. Computing the metric projection through B-spline approximation

The projected methods defined in Section 3 depend heavily on the availability of projection operators on the closed convex cone $X = \varphi_\mu(\mathcal{W}_2(\mathbb{R}))$. Being X a cone inside a linear space, such operators are always well defined, but their implementation might be nontrivial. In this Section, we present a possible solution to this problem, based on choosing a particular μ as base point and constructing a B-spline representation of the cone X .

4.1 Choosing μ as the uniform distribution on $[0, 1]$

As already mentioned, our projected methods can be carried out by choosing μ arbitrarily and there is no theoretical difference between different choices of μ , cf. Section 2.2. Nonetheless, in practice, a clever choice of μ can lead to substantially easier and more numerically stable algorithms. For instance, by choosing a measure μ with compact support C in \mathbb{R} , then the ambient space becomes $L_2^\mu(C)$ since we work up to zero-measure

sets. This greatly simplifies any numerical procedure since we could work with grids over bounded sets, and do not need to resort to any truncation procedure, which would be mandatory in case the support of μ was unbounded. Moreover, note that evaluating the maps φ_μ in a certain measure ν amounts to computing the transport map $T_\mu^\nu = F_\nu^- \circ F_\mu$, hence it is clear that the choice of F_μ numerically influences the results.

For the aforementioned reasons, we argue that a reasonable choice is to center our analysis in $\mu = U([0, 1])$. In fact, in this case, $L_2^\mu(\mathbb{R}) = L_2([0, 1])$, and $F_\mu = id_{[0, 1]}$ (the transport maps are simply given by quantile functions).

4.2 Metric Projection

Having chosen μ as Section 4.1 leads to an explicit characterization of the image of φ_μ as the set of square integrable a.e. non-decreasing functions on $[0, 1]$. Hence, the operator Π_X in Section 3 is the metric projection onto the cone of a.e. non-decreasing functions in $L_2([0, 1])$.

Projection onto monotone functions has been widely studied in the field of *order restricted* inference, (Anevski et al., 2006; Dykstra et al., 2012). For instance, in Anevski and Soulier (2011) an explicit characterization of such a projection is given, which however does not lead to a closed form solution, while in Ayer et al. (1955) several numerical algorithms to approximate the projection operator are proposed. Those algorithms are based on approximating the function to be projected with a step function defined on n intervals and can be shown to have a computational complexity that is linear in n (Best and Chakravarti, 1990).

Despite the numerical convenience of the aforementioned approximations, we believe that they are not suited for distributional data analysis. First and foremost, suppose that observations are given as probability density functions, so that one may want to interpret the results of a PCA, for instance, in terms of pdfs and not of quantile functions. If one were to estimate discontinuous principal directions through any of the algorithms in Ayer et al. (1955), it would not be possible to do so, as the corresponding cdfs would not be differentiable. In addition to that, the choice of the number of intervals n is not obvious when quantile functions are not directly observed but obtained with transformation. If n needs to be big to faithfully approximate the true quantile functions, this projection can be quite slow.

For these reasons, we propose to resort to a B-spline expansion, through which we can derive an alternative approximation of the projection operator Π_X , without incurring in the issues of the algorithms in Ayer et al. (1955). Moreover, we will also show in Section 5.3 that the proposed B-spline expansion also leads us to a simpler and faster reformulation of the geodesic PCA in Bigot et al. (2017).

4.3 Monotone B-splines representation

In what follows, let $\mu = U([0, 1])$. Moreover, denote with $\mathbf{x} = [x_1, \dots, x_k]' \in \mathbb{R}^k$ a generic vector.

As already said, through the φ_μ map, we can identify $\mathcal{W}_2(\mathbb{R})$ with the space

$$L_2([0, 1])^\uparrow := \{F^- \in L_2([0, 1]) \text{ s.t. } F^- \text{ is monotonically nondecreasing}\}$$

This leads us to consider a suitable B-spline basis for the space, to efficiently evaluate all the computations needed in our algorithms and for a convenient way to express the constraints which define $L_2([0, 1])^\uparrow$. In particular, we consider the basis of quadratic splines with equispaced knots in $[0, 1]$. The reason for this particular choice is two-folded. First

of all, splines of degree greater than one enjoy the nice property of uniform approximation of all continuous functions as the maximum distance between knots goes to zero, in turn this means that the closure of the linear space generated by the spline basis w.r.t the L_2 norm coincides with $L_2([0, 1])$. Secondly, quadratic splines are particularly well suited to characterize monotonic functions by looking at the coefficients of the (quadratic) B-spline expansion, as shown in the next Proposition.

Proposition 4 *Let $\{\psi_j^k\}_{j=1}^J$ be a basis of B-splines of order k defined over the knots x_1, \dots, x_{J+k+2} . Let $f(x) = \sum_{j=1}^J a_j \psi_j^k(x)$, then:*

1. *If the coefficients $\{a_j\}$ are monotonically increasing (decreasing) f is monotonically increasing (decreasing)*
2. *If $k = 2$, then 1. holds with an “if and only if”*

Before proceeding, let us fix some notation. From now on, we omit the dimension index “ k ” for the spline basis, writing ψ_j for ψ_j^2 , moreover we will let $\{\psi_j\}_{j=1}^J$ with fixed $J > 0$ denote a B-spline basis in $L_2([0, 1])$.

Remark 7 *Let $\mathbb{R}^{J\uparrow}$ be the set of vectors $v \in \mathbb{R}^J$ with nondecreasing coefficients. That is, letting $G = \{g_{ij}\}$ be the $J \times J$ binary matrix such that $\sum_j g_{ij} v_j = v_i - v_{i-1}$, for any element $v \in \mathbb{R}^J$ it holds that $Gv \geq 0$. Using Proposition 4, through the coordinates operator, the set $L_2([0, 1])^\uparrow \cap \text{Span}\{\psi_j\}_{j=1}^J$ is fully identifiable with $\mathbb{R}^{J\uparrow}$, endowed with the metric given by the symmetric positive definite matrix E with entries*

$$E_{ij} = \langle \psi_i, \psi_j \rangle_{L_2([0,1])}. \quad (9)$$

The norm induced is therefore $\|\mathbf{x}\|_E^2 = \mathbf{x}^T E \mathbf{x}$.

Remark 8 *It is possible to find a basis for \mathbb{R}^J with vectors lying in $\mathbb{R}^{J\uparrow}$ (and so in X_J), namely the vectors $(0, \dots, 0, 1)$, $(0, \dots, 0, 1, 1)$ etc. In other words, $\text{Span}(L_2([0, 1])^\uparrow \cap \text{Span}\{\psi_j\}_{j=1}^J) = \text{Span}\{\psi_j\}_{j=1}^J$ for every $J > 0$. This tells us that the convex cone of monotone splines is indeed quite big inside the spline space, and this a priori is beneficial for extrinsic methods, especially for PCA.*

From now on, to lighten the notation, we deliberately confuse the coefficients of the splines, living in \mathbb{R}^J or $\mathbb{R}^{J\uparrow}$ (with the metric given by E), with the corresponding spline functions living in the subsets of $L_2([0, 1])$ given by $L_2([0, 1])^\uparrow \cap \text{Span}\{\psi_j\}_{j=1}^J$ and $\text{Span}\{\psi_j\}_{j=1}^J$.

Remark 9 *Lastly, we point out that $\mathbb{R}^{J\uparrow}$ has the structure of a convex polytope, since the constraints given by $Gv \geq 0$ (guaranteeing that $v \in \mathbb{R}^{J\uparrow}$) are linear. Such geometric property makes optimization on $\mathbb{R}^{J\uparrow}$ handy and is key for the empirical methods developed in the remaining of the paper.*

As a consequence of Remark 9, the optimization problem given by the projection of a vector $v \in \mathbb{R}^J$ onto $\mathbb{R}^{J\uparrow}$ can be formulated as follows:

$$\Pi_{\mathbb{R}^{J\uparrow}}(v) = \arg \min_{Gw \geq 0} \|v - w\|_E. \quad (10)$$

The computational complexity required to solve (10) is at most cubic in the number of basis elements J (Potra and Wright, 2000).

Preliminary analysis showed that solving the optimization problem in (10) compares favorably with the Pool Adjacent Violators Algorithm (PAVA) in Ayer et al. (1955). In particular, computing PAVA with $n = 100$ approximation intervals is roughly eight times slower than (10) with $J = 20$ (a reasonable choice, leading to negligible approximation error, in our examples, with a quadratic spline basis). Increasing $n = 1000$ for PAVA makes it 700 times slower than (10).

In addition to that, resorting to a discretized approximation of quantiles would also increase the cost of the projected PCA, due to the need of using some functional PCA implementation, as opposed to the low-dimensional multivariate model we are able to implement with the B-spline basis functions.

5. Empirical Models with B-splines

In this Section, we present the empirical counterparts of the projected PCA defined in Section 3 and provide an illustrative example of projected linear regression, namely when both the dependent and independent variables are distributions.

Let $\{\psi_j\}_{j=1}^J$ be a fixed quadratic B-spline basis. Upon approximating the observed quantile functions with their spline expansion, thanks to Remark 7, we can develop our methodology in \mathbb{R}^J , considering the metric induced by E instead of the usual one. Indeed, given a vector $\mathbf{w} \in \mathbb{R}^J$, we can identify the corresponding function in L_2 by the map $\mathbf{w} \mapsto \sum_{j=1}^J w_j \psi_j$.

For the projected PCA in Section 5.1 and for the geodesic PCA in Section 5.3 we consider observations F_1^-, \dots, F_n^- , and let F_0^- be the centering point of the PCA. In our examples, F_0^- will always be the barycenter of the observations. As a preprocessing step, we approximate each of these quantile functions through a B-spline expansion and denote by $\mathbf{a}_i = \{a_{ij}\}_j$ and $\mathbf{a}_0 = \{a_{0j}\}_j$ the coefficients of the spline representation associated to F_i^- and F_0^- respectively, that is, $F_i^- \approx \sum_{j=1}^J a_{ij} \psi_j$. For the projected regression in Section 5.2, let observations $\{(F_z^-, F_y^-)_i\}_{i=1}^n$, where the F_{zi}^- 's are realizations of the independent variable \mathcal{Z} and the F_{yi}^- 's are realizations of the dependent variable \mathcal{Y} . We apply the same preprocessing step and let $\mathbf{a}_i^{(z)}$ and $\mathbf{a}_i^{(y)}$ denote the coefficient of the spline approximation of F_{zi}^- and F_{yi}^- respectively.

5.1 Empirical PCA

Denote with A the $(n \times J)$ matrix with rows $\mathbf{a}_1, \dots, \mathbf{a}_n$. As in standard PCA, the first principal component centered in \mathbf{a}_0 is found by solving the optimization problem:

$$\mathbf{w}_1^* = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_E=1} \sum_i |\langle \mathbf{a}_i - \mathbf{a}_0, \mathbf{w} \rangle_E|^2 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_E=1} \|A E \mathbf{w}\|^2 \quad (11)$$

where A is the matrix whose i -th row is given $\mathbf{a}_i - \mathbf{a}_0$. The optimization problem (11) can be solved similarly to a Rayleigh quotient: using Lagrange multipliers, (11) is equivalent to

$$\mathcal{L}(\mathbf{w}) := \mathbf{w}^T (A E)^T A E \mathbf{w} - \lambda (\mathbf{w}^T E \mathbf{w} - 1) \quad (12)$$

Deriving (12) w.r.t \mathbf{w} and equating the derivative to zero shows that the solutions to $d\mathcal{L}(\mathbf{w})/d\mathbf{w} = 0$ are the eigenvectors of the matrix $A^T A E$. Hence, ordering the eigenvalues of $A^T A E$ in decreasing order, the first principal component \mathbf{w}_1^* corresponds to the first eigenvector. Using similar arguments it can be shown that $\mathbf{w}_2^*, \dots, \mathbf{w}_J^*$ correspond to the remaining eigenvectors.

Once the first k principal directions $\mathbf{w}_1^*, \dots, \mathbf{w}_k^*$ are found, the projection of a new observation $x^* = \sum_{j=1}^J a_j^* \psi_j$ onto U_X^{k, x_0} (see Definition 1) is found exploiting Proposition 2. In particular, the following optimization problem is to be solved:

$$\begin{aligned} \arg \min_{\lambda_j \in \mathbb{R}} & \| (\langle \mathbf{a}^* - \mathbf{a}_0, \mathbf{w}_i^* \rangle_E - \lambda_i)_{i=1}^k \| \\ \text{s.t. } & G \left(\sum_{i=1}^k \lambda_i \mathbf{w}_i^* + \mathbf{a}_0 \right) \geq 0 \end{aligned} \quad (13)$$

which is equivalent to the minimization of a norm inside a polytope, that is a well-studied problem in \mathbb{R}^J (see Sekitani and Yamamoto, 1993) and there exist a variety of fast numerical routines to solve it.

5.2 Empirical Regression

In this section, we provide the details of the estimation procedure for a projected regression model where both the independent and the dependent variables are distribution-valued. It is straightforward to extend our methodology to cases when only one of these variables is distribution-valued and the other one takes values in \mathbb{R}^q .

First, we outline how to obtain an estimator for the linear operator Γ in Definition 4. Following Section 3.2 we first embed both \mathcal{Y} and \mathcal{Z} in $L_2([0, 1])$ through the inclusion operator $L_2([0, 1])^\uparrow \hookrightarrow L_2([0, 1])$, and assume the functional linear model presented in Ramsay (2004) and Prchal and Sarda (2007)

$$\mathcal{Y}(t) = \alpha(t) + \int_0^1 \beta(t, s) \mathcal{Z}(s) ds + \varepsilon(t), \quad t \in [0, 1] \quad (14)$$

so that $\Gamma = \Gamma_{\alpha, \beta}$ is the operator $\Gamma_{\alpha, \beta}(v)(t) = \alpha(t) + \int_0^1 \beta(t, s)v(s)ds$. The goal is then to estimate $\alpha \in L_2([0, 1])$ and $\beta \in L_2([0, 1]^2)$. Further, we assume that ε and \mathcal{Z} are uncorrelated: $\mathbb{E}[\mathcal{Z}(s)\varepsilon(t)] = 0$ for every $t, s \in [0, 1]$.

Consider now observations $\{(F_z^-, F_y^-)_i\}_{i=1}^n$ and the corresponding spline coefficients. Further, we project $\alpha(t)$ on the same spline basis, so that $\alpha \approx \sum_{j=1}^J \theta_{\alpha j} \psi(j)$ and $\beta(t, s)$ on the basis on $[0, 1]^2$ with $J \times J$ elements, so that $\beta(t, s) \approx \sum_{i, j=1}^J \Theta_{\beta ij} \psi_i(t) \psi_j(s)$. Neglecting the spline approximation error, model (14) entails

$$\mathbf{a}_i^{(y)} = \boldsymbol{\theta}_\alpha + \Theta_\beta E \mathbf{a}_i^{(z)} + \mathbf{a}_i^{(\varepsilon)}, \quad i = 1, \dots, n \quad (15)$$

where $\mathbf{a}_i^{(\varepsilon)}$ denotes the spline expansion coefficients of the unobserved error $\varepsilon_i(t)$.

We propose to estimate (15) using the same approach of Prchal and Sarda (2007), but extending it to account for spline approximations for both dependent and independent variables. We focus only on the estimate $\hat{\Theta}_\beta$ of Θ_β since once such estimate is obtained, the estimate for \mathbf{a}_α can be straightforwardly derived, (see Cai and Hall, 2006) as:

$$\hat{\boldsymbol{\theta}}_\alpha = \overline{\mathbf{a}^{(y)}} - \hat{\Theta}_\beta E \overline{\mathbf{a}^{(z)}}$$

where $\overline{\mathbf{a}^{(y)}}$ and $\overline{\mathbf{a}^{(z)}}$ are the means of $\mathbf{a}^{(y)}$ and $\mathbf{a}^{(z)}$ respectively.

The estimator $\hat{\Theta}_\beta$ is found by penalized least square minimization:

$$\hat{\Theta}_\beta = \arg \min_{\Theta} \frac{1}{n} \sum_{i=1}^n \left\| \left(\mathbf{a}_i^{(y)} - \overline{\mathbf{a}^{(y)}} \right) - \Theta E \left(\mathbf{a}_i^{(z)} - \overline{\mathbf{a}^{(z)}} \right) \right\|^2 + \rho \text{Pen}(1, \Theta) \quad (16)$$

where $\rho > 0$ is a penalization parameter to be fixed (usually through cross-validation) and $\text{Pen}(1, \Theta)$ is a penalization term defined in Prchal and Sarda (2007).

Briefly, the term $\text{Pen}(1, \Theta)$ in (16) penalizes both the norm of $\beta(t, s)$ and its derivatives, thus favoring smoother solutions. As shown in Prchal and Sarda (2007), (16) has a closed form solution. Nonetheless, the form of our solution differs from the one presented in Prchal and Sarda (2007), since they work directly on discretized functions while we propose to estimate spline coefficients, and some care must be taken since they can use (up to scaling) the usual inner product in the Euclidean space of discretized functions, while we must consider the inner product induced by E . However, the procedure for obtaining our result is identical to the one in Prchal and Sarda (2007). Hence, we only report the expression for the estimate.

Let \hat{C} be the matrix with entries

$$\hat{C}_{ks} = \left\langle \frac{1}{n} \sum_{i=1}^n \langle \mathbf{a}_i^{(z)}, b_k \rangle_E \mathbf{a}_i^{(z)}, b_s \right\rangle_E,$$

where b_k and b_s are the k -th and s -th elements of the standard Euclidean basis in \mathbb{R}^J . Further let \hat{D} the matrix with entries

$$\hat{D}_{ks} = \left\langle \frac{1}{n} \sum_{i=1}^n \langle \mathbf{a}_i^{(z)}, b_k \rangle_E \mathbf{a}_i^{(y)}, b_s \right\rangle_E.$$

Finally, let E' denote the matrix with entries $E'_{ij} = \langle \psi'_i, \psi'_j \rangle$ (where ψ'_i denotes the first derivative of the B-spline basis function ψ_i), $C_\rho = E^T \otimes (\hat{C} + \rho E')$, and $P = E'^T \otimes E + E^T \otimes E'$, where \otimes denotes the Kronecker product. Then the solution of (16) can be expressed as

$$\text{vec}(\hat{\Theta}_\beta) = (C_\rho + \rho P)^{-1} \text{vec}(\hat{D})$$

where $\text{vec}(\cdot)$ denotes the *vectorization* of the matrix.

Finally, our projected regression model is the composition of the operator induced by $(\hat{\theta}_\alpha, \hat{\Theta}_\beta)$ with the projection on $\mathbb{R}^{\uparrow J}$:

$$\mathbb{E}[\mathbf{a}_i^{(y)} \mid \mathbf{a}_i^{(z)}] = \Gamma_P(\mathbf{a}_i^{(z)}) = \Pi_{\mathbb{R}^{\uparrow J}} \left(\hat{\theta}_\alpha + \hat{\Theta}_\beta E \mathbf{a}_i^{(z)} \right).$$

5.3 An alternative optimization routine for the geodesic PCA and a comment on the computational costs

We now show how the framework in Section 4 can be employed also to derive faster numerical algorithms to find the global and nested geodesic PCA as of Definition 5 and Definition 6.

Proposition 5 (*Global geodesic PCA*) *A k dimensional global geodesic PC centered in \mathbf{a}_0 is the subset of $\mathbb{R}^{\uparrow J}$ spanned by $\{\mathbf{w}_1, \dots, \mathbf{w}_k\}$, linearly independent, which solve:*

$$\begin{aligned} \arg \min_{\{\lambda_i\}_1^n, \{\mathbf{w}_j\}_1^k} & \left\| \sum_{i=1}^n \mathbf{a}_i - \mathbf{a}_0 - \sum_{j=1}^k \lambda_{ij} \cdot \mathbf{w}_j \right\|_E^2 \\ \text{s.t. } & G \left(\sum_j \lambda_{ij} \mathbf{w}_j + \mathbf{a}_0 \right) \geq 0 \end{aligned} \tag{17}$$

Proposition 6 (*Nested geodesic PCA*) *With the same notation as above, a k dimensional nested geodesic PC, centered in \mathbf{a}_0 is the set spanned by $\{\mathbf{w}_1, \dots, \mathbf{w}_k\}$ in $\mathbb{R}^{J\uparrow}$, where the \mathbf{w}_i s are found recursively from \mathbf{w}_1 to \mathbf{w}_k , such that \mathbf{w}_h is a solution, for every h , of:*

$$\begin{aligned} \arg \min_{\{\lambda_i\}_{i=1}^n, \mathbf{w}} \sum_{i=1}^n \|\mathbf{a}_i - \mathbf{a}_0 - \lambda_i \mathbf{w}\|_E^2 \\ \text{s.t. } \langle \mathbf{w}_j, \mathbf{w} \rangle_E = 0, \quad j = 1, \dots, h-1 \\ G(\lambda_i \mathbf{w} + \mathbf{a}_0) \geq 0, \quad \|\mathbf{w}\|_E = 1 \end{aligned} \quad (18)$$

To solve (17) and (18) we employ an interior point method using the solver Ipopt (Wächter and Biegler, 2006). When comparing our implementation with $J = 20$ spline basis and the one in Cazelles et al. (2018), we notice a substantial performance improvement, by a factor of 35 for a data set of $n = 100$ distributions, due to the fact working with spline approximations reduces greatly the number of parameters in the optimization problem.

Further, note that (17) and (8) seem extremely similar. However, in (8) the optimization is carried out having fixed $\mathbf{w}_1^*, \dots, \mathbf{w}_k^*$ and for a single observation, while in (17) the optimization is done over a much larger set of parameters. In fact, the number of parameters in (17) is $(n+k)J$, hence the computational complexity needed to solve (17) is cubic in both the number of bases and the number of observations. On the other hand, the projected PCA requires a linear time in the number of observations (computation of $A^T A E$) and cubic time in the number of basis J (eigendecomposition and projections of new observations).

6. Asymptotic Properties

In this section, we study the convergence of the proposed projected empirical methods. First of all, we show that as the number of spline basis J increases, the error due to the spline approximation vanishes if the data is sufficiently regular. Further, under a suitable set of assumptions, we establish consistency results for the projected PCA and for the projected distribution on distribution regression.

6.1 Convergence of Quadratic B-splines

In the following, denote with $W_k^r([0, 1])$ the space of functions whose weak derivatives up to order k belong to $L_r([0, 1])$, further denote with D the (weak) derivative operator, so that $Df = f'$, $D^2f = f''$ and so on,

Proposition 7 *Let μ a probability measure on \mathbb{R} , F_μ^- its quantile function such that $F_\mu^- \in W_3^\infty$. For each J let $\{\psi_j\}_{j=1}^J$ denote a quadratic B-spline basis on J equispaced knots in $[0, 1]$. Then there exist a sequence of spline functions $S_J = \sum_{j=1}^J \lambda^{(j)} \psi_j^{(j)}$, with $\lambda_j^{(j)}$ monotonically non-decreasing in j for every J , such that:*

$$\|S_J - F_\mu^-\|_\infty \leq C \|D^2 f_\mu^-\|_\infty J^{-2}$$

with $f_\mu^- = DF_\mu^-$ and $C > 0$ constant.

Let us remark two important facts.

Remark 10 Since the inclusion $L_\infty([0, 1]) \subset L_2([0, 1])$ is continuous, thanks to Hölder inequality, the convergence rates hold also for the L_2 norm. By default we will use the L_2 norm if not stated differently.

Remark 11 By Poincaré inequality, if $\|D^3 f\|_\infty < C$ then f belongs to a sphere in $W_3^\infty([0, 1])$ whose radius depends on C and on the Poincaré constant of $[0, 1]$; viceversa, all the elements in the sphere of radius C in $W_3^\infty([0, 1])$ clearly have (weak) derivatives bounded by C .

6.2 Consistency

In this Section we prove the consistency of the projected methods under some assumptions on the data-generating process. In particular, we show that there exists a number of basis functions $J > 0$ and a sample size n such that the error committed by the empirical models in Section 5 is smaller than $\varepsilon > 0$, for any fixed ε .

6.2.1 PCA

Consistency of spline-based PCA for functional data has been addressed, among the first, by Silverman et al. (1996) and Qi and Zhao (2011). As one of the main building blocks of our projected PCA is the PCA in the ambient space, that is $L_2([0, 1])$, it is natural to follow Qi and Zhao (2011) in making the following assumptions. Consider data $\mu_1, \dots, \mu_n, F_1^-, \dots, F_n^-$ the corresponding quantile functions, then:

- (P1) The data generating process satisfies $F_1^-, \dots, F_n^- \sim \mathcal{F}$ with the F_i^- independent and $\mathbb{E}[\mathcal{F}] = 0$.
- (P2) F_1^-, \dots, F_n^- can be approximated by functions in W_3^∞ with uniformly bounded third derivative.
- (P3) $\mathbb{E}[\|F_i^-(t)\|^4] < \infty, i = 1, \dots, n$.
- (P4) The eigenvalues of the covariance operator of \mathcal{F} have multiplicity 1.
- (P5) The eigenfunctions of the covariance operator of \mathcal{F} belong to some bounded set in $W_3^\infty([0, 1]) \subset W_3^2([0, 1])$.

Before stating the main results, let us comment on assumptions (P1)-(P5). First of all, (P2) is essential in order to apply Proposition 7 and get uniform errors on the data set. Moreover, (P2) is satisfied, for instance, if the F_i^- 's lie in the L_2 -closure of a ball of radius $M > 0$ in W_3^∞ . (P4) is a rather standard condition and is satisfied if $\mu_1, \dots, \mu_n \in \mathcal{W}_4(\mathbb{R})$. (P4) and (P5) imply the assumptions that in Qi and Zhao (2011) are used for the consistency results. In particular, (P5) is stronger than the corresponding assumption in Qi and Zhao (2011), where the eigenfunctions are assumed to belong to $W_2^2([a, b])$. Similarly, in such work, there is no counterpart of assumption (P2); in fact we need these stronger regularity conditions to get uniform errors when using B-splines. Still some of the examples Qi and Zhao (2011) provide of situations satisfying their assumptions, meet also our requirements. Finally, the zero-mean assumption in (P1) might seem a little odd, since we know that the quantile functions are monotonically nondecreasing. However, observe that it is always possible to subtract the empirical mean from the observations to satisfy (asymptotically) this assumption.

Let J denote the dimension of a quadratic B-spline basis on $[0, 1]$ and let \mathbf{a}_i^J the coefficients of the B-spline approximation of F_i^- . In what follows, to lighten the notation,

we refer to a set of spline coefficients both as elements of \mathbb{R}^J with the E -norm, or as functions in L_2 , without making explicit reference to the coordinate operator and its inverse.

Proposition 8 *Under assumptions (P1)-(P5), for any $\varepsilon > 0$ there exists a sample size $n > 0$ and a number of basis functions $J > 0$ such that:*

$$\left| \max_{\|w\|_{L_2}=1} \frac{1}{n} \sum_i \langle F_i^-, w \rangle_{L_2}^2 - \max_{\|w\|_E=1} \frac{1}{n} \sum_i \langle \mathbf{a}_i^J, w \rangle_E^2 \right| < K\varepsilon$$

for some constant $K > 0$.

Proposition 8 ensures the consistency of the B-spline approximation of the PCA for mono-functional data in H which is equivalent to the consistent estimation of the projected principal directions.

Suppose now to have computed $U_k^J = \{\mathbf{w}_h^{J*}\}_{h=1}^k$, that is the approximations of the principal directions $U_k = \{\mathbf{w}_h^*\}_{h=1}^k$ found with J basis functions. We observe that $Sp(U_k^J) \cap L_2([0, 1])^\uparrow = Sp(U_k) \cap \mathbb{R}^{J\uparrow}$. Since for any set of coefficients λ_h we have the convergence $\sum \lambda_h \mathbf{w}_h^{J*} \rightarrow \sum \lambda_h \mathbf{w}_h^*$, we obtain that the projection of a point onto $Sp(U_k^J) \cap L_2([0, 1])^\uparrow$ converges to the projection onto $Sp(U_k) \cap L_2([0, 1])^\uparrow$. Thus we also have convergence of the projection onto the principal components.

6.2.2 REGRESSION

We consider model (14) given samples $\{(F_z^-, F_y^-)_i\}_{i=1}^n$. We make the following assumptions:

(R1) The data generating process satisfies (14) and $\mathbb{E}[\mathcal{Z}(s)\varepsilon(t)] = 0$ for every $t, s \in [0, 1]$.

(R2) $\alpha \in L_2([0, 1])$ and $\beta \in L_2([0, 1] \times [0, 1])$.

(R3) With probability 1, each quantile function in the samples $\{(F_z^-, F_y^-)_i\}_{i=1}^n$ lies inside a sphere of radius $K > 0$ in $W_\infty^3([0, 1])$.

Without loss of generality, suppose that both the dependent and the independent variables have been centered by subtracting their mean so that $\mathbb{E}[\mathcal{Z}] = \mathbb{E}[\mathcal{Y}] = 0$ and $\alpha = 0$.

The strategy to prove the consistency of the projected linear regression is the following. First of all, we prove that the estimator $\hat{\Theta}_J$ converges to the estimator $\hat{\Theta}_{PS}$, defined in Prchal and Sarda (2007), for large enough n and J . Second, we exploit the consistency of the estimator in Prchal and Sarda (2007) combined with the approximation results of the metric projection, to establish consistency in terms of the prediction error of our projected regression operator.

Briefly $\hat{\Theta}_{PS}$ is obtained by minimizing an objective function similar to the one in (16), but where the spline approximation is used only for Θ , while the F_{zi}^- 's and the F_{yi}^- 's are assumed fully observed, and not approximated through splines. Calling B the vector of functions with entries ψ_1, \dots, ψ_J , $\hat{\Theta}_{PS}$ is defined as:

$$\hat{\Theta}_{PS} = \arg \min_{\Theta} \frac{1}{n} \sum_i \|F_{yi}^- - \langle F_{zi}^-, B^T \Theta B \rangle\|^2 + \rho \text{Pen}(1, \Theta).$$

Convergence of $\hat{\Theta}_J$ to $\hat{\Theta}_{PS}$ is shown in the next proposition

Proposition 9 *Under assumptions (R1)-(R3), if the number of samples is big enough $\widehat{\Theta}$ and $\widehat{\Theta}_J$ exists with probability close to 1, and there is $J > 0$ such that $\|\widehat{\Theta}_{PS} - \widehat{\Theta}_J\|_{E \otimes E} < \varepsilon$.*

Let $\widehat{\beta}_{PS}$ and $\widehat{\beta}_J$ be the kernels $\widehat{\beta}_{PS} = B^T \widehat{\Theta}_{PS} B$ and $\widehat{\beta}_J = B^T \widehat{\Theta}_J B$. Since $\|\widehat{\beta}_{PS}(s, t) - \widehat{\beta}_J(s, t)\|_{L_2([0,1]^2)} = \|\widehat{\Theta}_{PS} - \widehat{\Theta}_J\|_{E \otimes E}$, we established strong convergence of our kernel to the estimator of Prchal and Sarda (2007). This implies that the consistency results for the estimator $\widehat{\Theta}_{PS}$ holds also for $\widehat{\Theta}_J$, with respect to the seminorm induced by the covariance operator of \mathcal{Z} .

Specifically, given \mathcal{Z} H -valued random variable and its covariance operator $\mathcal{C}_{\mathcal{Z}}$, for any $\varphi \in L_2([0, 1]^2)$, we consider the semi-norm on $L_2([0, 1]^2)$ given by:

$$\|\varphi\|_{\Gamma_{\mathcal{Z}}} = \int_{[0,1]} \langle \mathcal{C}_{\mathcal{Z}} \varphi(\cdot, t), \varphi(\cdot, t) \rangle dt$$

Thus, the following result is immediately implied since strong convergence implies seminorm convergence (see Appendix A).

Corollary 1 *For $J > 0$ big enough $\mathbb{E}[\|\beta - \widehat{\beta}_J\|_{\mathcal{C}_{\mathcal{Z}}}] < \varepsilon$.*

Proof *We use the seminorm triangle inequality:*

$$\|\beta - \widehat{\beta}_J\|_{\mathcal{C}_{\mathcal{Z}}} \leq \|\beta - \widehat{\beta}\|_{\mathcal{C}_{\mathcal{Z}}} + \|\widehat{\beta} - \widehat{\beta}_J\|_{\mathcal{C}_{\mathcal{Z}}}.$$

The first term on the right hand side converges to zero thanks to Theorem 2 in Prchal and Sarda (2007), while the second term converges to zero thanks to Proposition 9 and the previous observations. ■

Lastly, we need to take into account the projection step. First, we notice that $\|\beta - \widehat{\beta}\|_{\Gamma_{\mathcal{Z}}}$ corresponds to the expected prediction error, in fact, as in Prchal and Sarda (2007):

$$\|\beta - \widehat{\beta}_J\|_{\mathcal{C}_{\mathcal{Z}}} = \int_{[0,1]} \mathbb{E} \left[\langle \mathcal{Z}, \beta(\cdot, t) - \widehat{\beta}_J(\cdot, t) \rangle^2 \mid \widehat{\beta}_J \right] dt,$$

further, by Hölder's inequality $\mathbb{E} \left[|\langle \mathcal{Z}, \beta - \widehat{\beta}_J \rangle| \mid \widehat{\beta}_J \right] \rightarrow 0$, which straightforwardly yields $\mathbb{E} \left[\|\Gamma_{\beta}(z) - \Gamma_{\widehat{\beta}_J}(z)\| \mid \widehat{\beta}_J \right] \rightarrow 0$.

Thus, the following simple lemma ensures the consistency of the spline approximation of the projection on X and leads to the consistency of the projected regression in terms of prediction error. Again, following Remark 7, we can identify the space monotone B -splines with J basis functions with $\mathbb{R}^{J \uparrow}$. Hence, to lighten the notation, we denote $\Pi_{\mathbb{R}^{J \uparrow}}$ the metric projection operator onto the space of monotone B -splines with J basis functions.

Lemma 2 *Given $\beta_n \rightarrow \beta$ in H , for any $\varepsilon > 0$ there exists $n, J > 0$ such that $\|\Pi_{\mathbb{R}^{J \uparrow}}(\beta_n) - \Pi_{L_2([0,1]) \uparrow}(\beta)\| \leq \varepsilon$.*

7. Numerical Illustrations for the PCA

In this section we perform PCA on different simulated data sets and on a real data set of Covid-19 mortality data in the US. In particular, on the simulated data sets we compare the performance of our projected PCA (in terms of approximation error and interpretability of the directions) with the ones of intrinsic methods, showing that the projected PCA is

a valid competitor in a diverse set of situations. For the Covid-19 data set, we compare inference obtained using the projected, nested and log PCA, highlighting the practical benefits of the projected PCA over the log one.

For the projected, nested and global PCAs we need to fix a B-spline basis to express the quantile functions. In particular, we fix an equispaced quadratic B-spline basis with J interior knots on $[0, 1]$. Here, the number of basis J is always fixed to 20, which provided a negligible approximation error of the quantile functions. We did not observe any appreciable change when increasing it. In Appendix C we show further simulations where we perform sensitivity analysis as the number of basis increases for a fixed sample size, we provide empirical confirmation of the consistency results in Section 6 and give practical guidance on how to choose J .

7.1 Simulation studies

We consider three different simulations to compare both the interpretability and the ability to compress information of different PCAs.

We compare our projected PCA with the nested and global geodesic PCAs (Bigot et al., 2017; Cazelles et al., 2018) and the *simplicial* PCA (Hron et al., 2014).

Briefly, the simplicial PCA applies a transformation that maps densities defined on the same compact interval I into functions in $L_2(I)$, called *centered log ratio*. Then, a standard L_2 PCA is performed on the transformed pdfs and, by the inverse of the centered log ratio transform, the results are mapped back to the space of densities, called Bayes space (for a more accurate definition, see Egozcue et al., 2006). In particular, we remark that, to be well defined, the simplicial PCA requires that all the pdfs have support equal to I , which is a strong assumption in practice. Further details about simplicial PCA are given in Appendix B.

As for the projected PCA, to compute the simplicial PCA, we resort to a B-spline approximation, but this time of the transformed pdfs. Hence, we need to select a B-spline basis on the support of the pdfs I . In this case, we fix a cubic B-spline basis with

$$J' = J = 20$$

interior knots on I , as this choice yielded a negligible approximation error for the transformed pdfs.

In the first scenario, we simulate data from

$$\begin{aligned} p_i(x) &\propto \frac{1}{\sigma_i} \exp\left(-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right) \mathbb{I}(x \in [-10, 10]), \quad i = 1, \dots, 100 \\ \mu_i &\sim 0.5\mathcal{N}(-3, (0.2)^2) + 0.5\mathcal{N}(3, (0.2)^2) \\ \sigma_i &\sim \text{Uniform}([0.5, 2.0]) \end{aligned} \tag{19}$$

Where “proportional to” stands for the fact that we confine the density to the support $[-10, 10]$ and renormalize it so that it integrates to 1.

Observe that there are two sources of variability across the pdfs from the data generating process (19). The first one is the location of the *peak* μ_i and the second one is the *width* of the distribution around the peak, controlled by σ_i . See Figure 3.

Figure 4 shows the first two principal directions obtained using the different methods. We can notice several differences between them. Focusing on the first principal direction, we can see that the simplicial, projected and nested PCAs detect a change in the location of the peak of the pdf. In particular, the first direction for the Wasserstein PCAs represents a shift from left to right of this peak, while for the simplicial PCA the first direction is

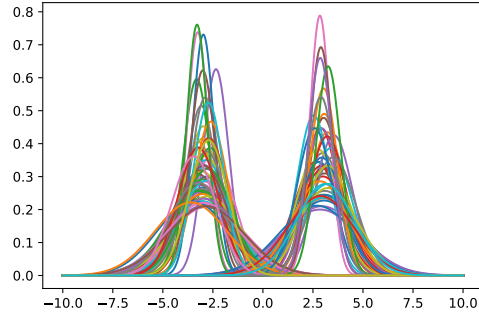


Figure 3: Data set of pdfs generated from (19)

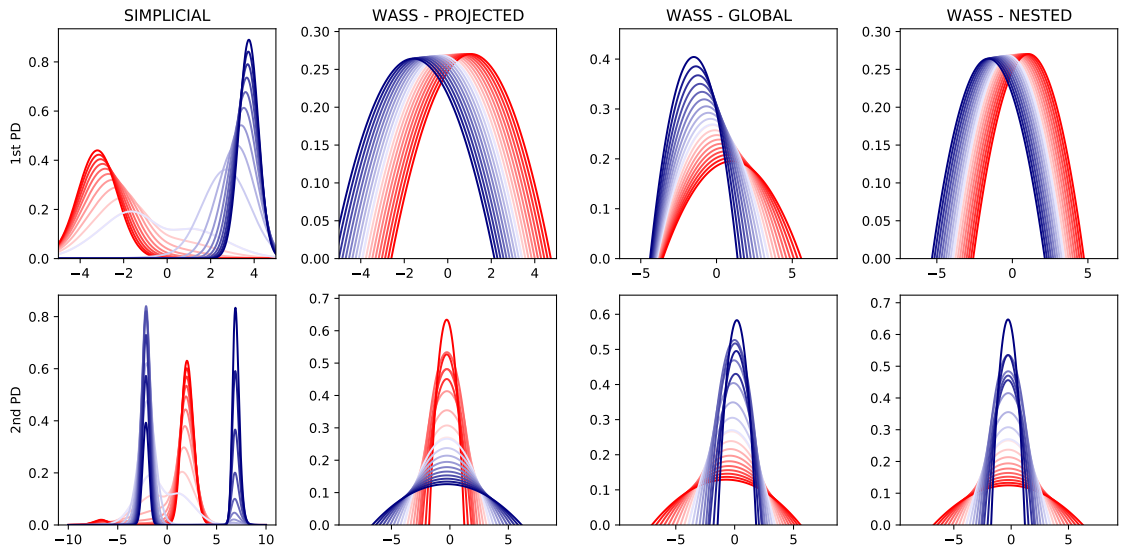


Figure 4: Top row: first principal direction. Bottom row: second principal direction. Each line represents the pdf associated to $\lambda \mathbf{w}_i$ where \mathbf{w}_i is the i -th principal direction ($i = 1, 2$) and λ is a score ranging from -2 (darkest blue) to $+2$ (darkest red).

associated to a peak in 3 (blue lines, negative values of the scores) or to a peak in -3 (red lines, positive value of the scores). This also highlights the difference in the geometries underlying the Wasserstein and Bayes spaces. Looking at the second principal direction instead, we can see how in the Wasserstein PCAs it clearly represents a change in the width of the distribution, while for the simplicial PCA the interpretation is somewhat obscure.

The global geodesic PCA deserves a separate discussion. Indeed, from Definition 5 it is clear that a global principal component is a convex set without any notion of preferential directions, so that it is not possible to interpret separately the variation along the first and second direction found by the global PCA.

Now we present two additional simulations that quantify the amount of information that is “lost” by performing the PCA. As a metric, we consider the reconstruction error, that is, the quantity

$$RE_k = \frac{1}{n} \sum_{i=1}^n W_2(F_i^-, \tilde{F}_i^-) \quad (20)$$

where the F_i^- ’s are the observed probability measures, \tilde{F}_i^- are the reconstructed ones and k is the dimension of the principal component. More in detail \tilde{F}_i^- is found by first projecting $(F_i^- - F_0^-)$ into \mathbb{R}^k using the PCA and then applying the inverse transformation. Informally, the reconstruction error is a measure of the quantity of information lost by applying the PCA as a black-box dimensionality reduction.

As evident in Equation (20), we measure the performance of PCAs just in terms of Wasserstein metric. This is likely to favor the performance of the Wasserstein PCAs over the simplicial one. Thus, the interesting performance comparison is the one between the geodesic PCAs and the projected PCA. Nevertheless, we think that is worth reporting also the results for the simplicial PCA, which is an intrinsic method in the Bayes space, to show that the underlying metric structures are extremely different. This also helps to appreciate the results in Section 8. Given the difference in the metric structure between Wasserstein and Bayes spaces, we believe that the choice between simplicial and Wasserstein frameworks is not trivial and should be application-driven.

To measure raw performance differences between geodesic and projected PCAs, we simulate data so that there is little recognizable structure in them, unlike in the previous example. The data generating process is as follows:

$$\begin{aligned} p_i(x) &\propto \sum_{j=1}^K w_{ij} \frac{1}{\sigma_{ij}} \exp\left(-\frac{(x - \mu_{ij})^2}{2\sigma_{ij}^2}\right) \mathbb{I}(x \in [-10, 10]) + 10^{-5}, \quad i = 1, \dots, 100 \\ \mathbf{w}_i &\sim \text{Dirichlet}_K(1/K) \\ (\mu_{ij}, \sigma_{ij}) &\sim \mathcal{N}(d\mu_{ij}; 0, 2^2) \text{Uniform}(d\sigma_{ij}, 0.5, 2.0) \end{aligned} \quad (21)$$

Observe that (21) is a finite dimensional approximation of the Dirichlet Process mixture model, a popular workhorse in Bayesian nonparametric statistics, that is well known to be dense in the space of densities on \mathbb{R} , see for instance Ferguson (1983). An example of the kind of pdfs generated from (21) is shown in Figure 5(a).

To separate the effect of the B-spline smoothing procedure, in this scenario we evaluate the reconstruction error in (20) considering $\tilde{\mu}_i$ to be the reconstructed quantile functions (for the Wasserstein PCAs) or pdfs (for the simplicial PCA) and μ_i to be the probability measure represented by the B-spline approximation of the quantile function or the (centered log ratio of) the pdf respectively.

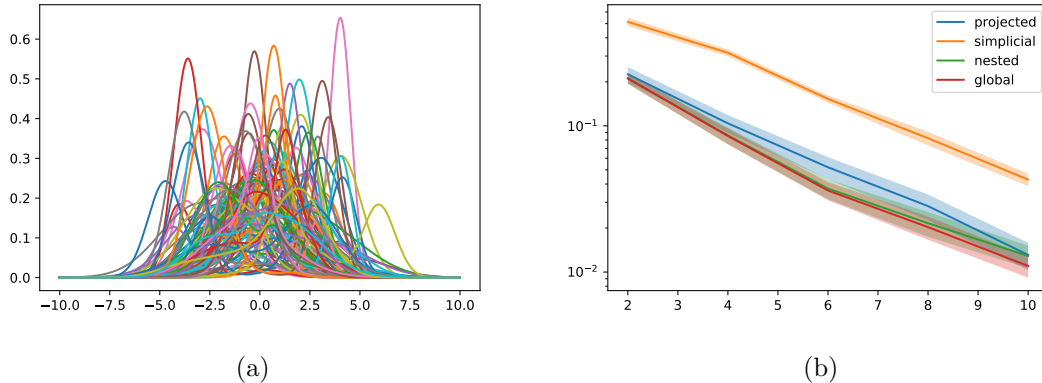


Figure 5: Left panel: example of simulated data set for Scenario 2. Right panel: reconstruction error as a function of the dimension of the principal component employed for the different methods. The solid lines represent the mean of 10 independent runs on independent data sets from (21) and the shaded area represent \pm one standard deviation.

Figure 5(b) shows the reconstruction error as a function of the dimension of the principal component, that is, RE_k as a function of k . We can see how the three Wasserstein PCAs consistently outperform the simplicial one. Moreover, as to be expected, the global geodesic PCA obtains the lowest reconstruction error for all the choices of dimension k , with the nested geodesic PCA being a close runner-up. However, the computational cost of finding the nested or global geodesic PCA can become prohibitive as the sample size or the number of bases in the B-spline expansion or the dimension k increases. For comparison, finding the 10-dimensional projected PCA is around 1,000 times quicker than finding the corresponding global geodesic PCA and 200 times quicker than finding the nested geodesic one.

As an additional simulation, in Appendix C we investigate the effect of the number of B-spline basis J . In particular, we conclude that, for a fixed dimension k the reconstruction error (20) increases with the number of basis functions, both for the projected and the simplicial PCA. Furthermore, we also observe that the reconstruction error for the simplicial PCA exhibits a larger variance than the reconstruction error for the projected PCA. Our insight is that this is due to the different degree of smoothness of the pdfs and the quantile functions. Since the quantile functions are in general smoother than the pdfs, their B-spline expansion should have lower variance.

7.2 Assessing the reliability of the projected PCA

A classical measure of performance of the standard Euclidean PCA, also useful to determine the dimension of the principal component to use, is the proportion of the explained variance. For a k -dimensional Euclidean principal component, this quantity is easily computed as a ratio of eigenvalues: $\sum_{j=1}^k \lambda_j / \sum_{j \geq 1} \lambda_j$. Upon truncating the series at the denominator, the same quantity can also be computed for PCA in infinite dimensional Hilbert spaces.

Due to the projection step involved in our definition of PCA, we argue that the proportion of explained variance might not be a reliable indicator of performance, nor should it be used to guide the choice of the dimension k . Instead, we propose a fast alternative based on the Wasserstein distance that we believe better represents the properties of the

projected PCA, that is, the normalized reconstruction error:

$$NRE_k = \frac{\frac{1}{n} \sum_{i=1}^n W_2(F_i^-, \tilde{F}_i^-)}{\frac{1}{n} \sum_{i=1}^n W_2(F_i^-, F_0^-)},$$

where the numerator corresponds to the reconstruction error in (20) and the denominator is the average distance between the observed measures and their barycenter. Observe that in Euclidean spaces, this quantity is closely related to the proportion of explained variance, since in Euclidean spaces maximizing variance in a subspace, amounts to minimizing the average distance from the subspace to data points.

Given its extrinsic nature, for a fixed dimension, the projected PCA might sometimes fail to capture the variability of some particular data set and, in those situations, an intrinsic approach should be preferred. However, given the high computational cost associated to geodesic PCAs, one would carry out such analysis only knowing that the results would be significantly better than the ones obtained by projected PCA. This calls for discerning whether the poor performance of projected PCA is due to its extrinsic nature or rather to the scarceness of structure in the data set under consideration: in the former situation it is likely that a geodesic approach would yield better results, in the latter instead, it is likely that results remain the same.

We propose now two empirical indicators of the “reliability” of the empirical projected PCA. The first one measures, once a k -dimensional principal component is found, how reliable are the projected principal directions and the second one gives an idea of how different the projected PCA and the L_2 PCA are. To assess the interpretability of the principal directions and the scores obtained with the projected PCA, we first compute for every principal direction \mathbf{w}_h^* the quantities η_h^{\min} and η_h^{\max} such that

$$\eta_h^{\min} = \min_{\eta \in \mathbb{R}} \{\mathbf{a}_0 + \eta \mathbf{w}_h^* \in \mathbb{R}^{J\uparrow}\}$$

where \mathbf{a}_0 is the spline coefficient vector associated with the barycenter F_0^- . The scalar η_h^{\max} is found analogously. Hence $(\eta_h^{\min} \mathbf{w}_h^*, \eta_h^{\max} \mathbf{w}_h^*)$ is the segment spanned by the principal direction living inside the convex cone $\mathbb{R}^{J\uparrow}$. If the scores of all observations along this direction lie within the range $(\eta_h^{\min}, \eta_h^{\max})$, then the variability captured by (empirical) projected PCA can be decomposed along the principal directions, whose scores are then highly interpretable. Contrary, the PCA scores outside $(\eta_h^{\min}, \eta_h^{\max})$ will be associated with functions which are not quantiles, and thus limiting the interpretability of the direction. Hence, we propose the following *interpretability score*

$$IS_h = 1 - \frac{1}{n} \sum_{i=1}^n d(s_{ih}, [\eta_h^{\min}, \eta_h^{\max}]) / |s_{ih}|, \quad (22)$$

where s_{ih} is the score of observation i along direction h according to the projected PCA. A value of IS_h equal to one corresponds to perfect interpretability, that is, projected PCA behaves like a standard Euclidean PCA along direction h . On the other hand, values of IS_h closer to zero indicate that the decomposition of the variance along the principal directions lies outside $\mathbb{R}^{J\uparrow}$ for direction h . The interpretability score can be fruitfully used also to evaluate the directions found with the nested PCA, upon replacing the s_{ih} 's in (22) with the scores given by the nested PCA.

Note that the IS_h score is useful to interpret the directions one at a time. However, it can be the case that some scores along one direction h' lie outside the $(\eta_{h'}^{\min}, \eta_{h'}^{\max})$ range but that the L_2 projection on the $h \geq h'$ component still lies within the projected component.

For instance, this could imply that a projected PC could be similar to a nested one despite having very different directions. A discrepancy between the two can appear when the projections of some data points on the L_2 PCA lie outside $\mathbb{R}^{J\uparrow}$. Using the terminology of Proposition 2 this can be measured in terms of difference between the projections $\Pi_k(F^{-*} - F_0^-)$ and $\Pi_{Sp(U_k) \cap (X-x_0)}(F^{-*} - F_0^-) = \Pi_{Sp(U_k) \cap (X-x_0)}(\Pi_k(F^{-*} - F_0^-))$, for a given observation F^{-*} . To quantify the loss of information at the level of the component (instead of direction), we propose to measure the “ghost variance” captured by the L_2 PCA:

$$GV_k = \frac{1}{n} \sum_{i=1}^n \|\Pi_k(F_i^- - F_0^-) - \Pi_{U_{X^-,k}}(F_i^- - F_0^-)\|_2 / \|F_i^- - F_0^-\|_2,$$

that is, the GV_k score measures the quantity of information that is lost due to the projection step or, in other words, the information that we trained our PCA on, but that does not appear in the Wasserstein Space. If $GV_k = 0$ then all the information captured by the L_2 PCA is inside the Wasserstein Space, then the projected PCA coincide with the nested one by definition.

Finally, although this situation never occurred in our experience, it might happen that GV_k is small but some IS'_k ($k' \leq k$) is large. This means that the subspace identified by the projected PCA is suitable for representing the data, but the single principal directions are not interpretable. In this case, we suggest to take a hybrid approach: use the projected PCA as a fast black-box dimensionality reduction step, thus reducing the dimensionality of each observation from J to k , and then use the nested PCA, in dimension k , to estimate the directions, the main advantage being the reduction in the computational cost to estimate the nested PCA in this lower dimensional space.

7.3 Analysis of the Covid-19 mortality data set

We perform PCA analysis on the Covid-19 mortality data publicly available at data.cdc.gov as of the first December 2020. The data set collects the total number of deaths due to Covid 19 in the US from January 1st 2020 to the current date, data are subdivided by state, sex, and age. In particular, the ages of the deceased are grouped in eleven bins: $[0, 1)$, $[1, 5)$, $[5, 15)$, $[15, 25)$, $[25, 35)$, $[35, 45)$, $[45, 55)$, $[55, 65)$, $[75, 85)$, $[85, +\infty)$ but we truncate the last bin to 95 years for numerical convenience. Further, we remove Puerto Rico from the analysis because it presented too many missing values. Our final data set, shown in Figure 6(a), consists of 106 samples of the distribution of the ages of patients deceased due to Covid-19, divided by sex and pertaining 53 between US states and inhabited territories.

We apply our usual B-spline approximation with $J = 20$ basis to the quantile functions obtained starting from the histograms in Figure 6. This choice of J yields an average approximation error, in terms of Wasserstein distance, of 0.02. An error this low is to be expected since the quantile functions are piecewise linear functions defined on eleven intervals.

We use this real data set to make a hands-on comparison of the inference that can be obtained employing the projected, nested and log PCA.

We start by comparing the projected and nested PCAs. The first direction found by the nested PCA is identical to the one found by the projected while the second is extremely close: the cosine between the two principal directions is approximately 0.99. In line with this, the interpretability scores equal $IS_1 = 1$ and $IS_2 \approx 0.89$, while $GV_2 = 0.05$. Moreover, the two-dimensional projected principal component explains more than 90% of the L_2

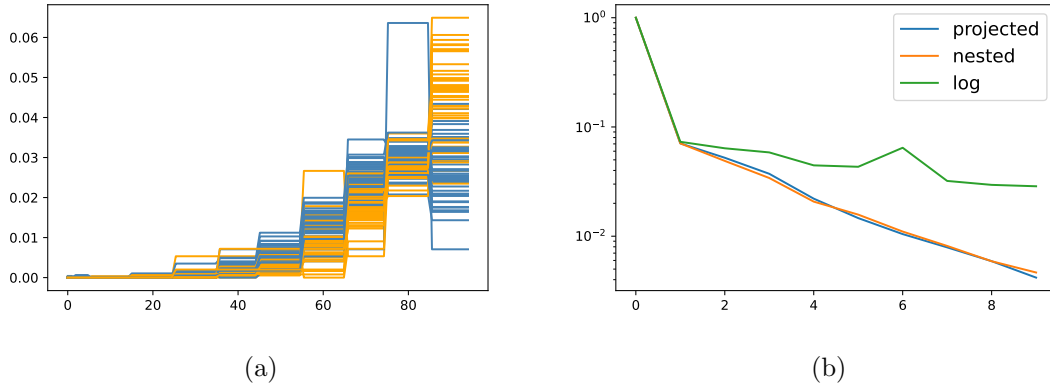


Figure 6: Left panel: distributions of age at the time of death for Covid-19 patients divided by sex: orange corresponds to female and blue to males. Different lines correspond to different US states / inhabited territories. Right panel: reconstruction error as a function of the dimension of the component for different PCAs. The 0-th principal component is the empirical mean.

variability and $NRE_2 \approx 0.05$ for both projected and nested PCA. Given the reconstruction error and the GV_2 score, we can conclude that the two-dimensional projected principal component provides a very good fit to the data, and that both selected principal directions are well behaved with respect to their scores, guaranteeing interpretable results.

Considering the discussion above and the fact that both the projected and nested PCA employ metric projection to map data points to the k -dimensional principal component, inference obtained with the nested PCA and with the projected one is almost identical in this case and we show results only for the projected PCA in Figure 7. In particular, the first principal direction shows that the greatest variability is due to the elders: low negative values along this direction correspond to most of the mortality being concentrated among in the 80+ range. The red and the green distributions shown in the rightmost panel show two antithetic behaviors which correspond to scores along the first principal direction of roughly -8.5 and 7 as shown in the third panel of Figure 7. In fact, the red distribution is concentrated almost exclusively on the last two bins of the histogram, with the 85+ bin weighting for more of 60% of the deaths. At the opposite, the green distribution gives more weight to lower age values. The second direction instead shows variability in the 40 – 80 range. The purple distribution, characterized by the highest score along this direction, shows that a significant percentage of deaths occurred in the age range 60 – 75. Finally, the third panel of Figure 7 reports the scores along the first two principal directions for the whole data set, blue dots representing males and orange dots women. We can appreciate how women tend to have lower scores on both directions. This is in line with our understanding that Covid-19 is more severe among the male population (see for instance Mandavilli, 2020), which explains why males are more susceptible to death even at younger ages, while deaths among women are more concentrated in the 70+ age range, being the elders in general more fragile.

The comparison with log PCA requires more attention. First of all, note that the directions obtained with the projected and log PCA are the same by definition, since they are both obtained performing PCA in $L_2([0, 1])$, but the principal components may differ because different projection operators are employed when the orthogonal projection of a point onto the principal component lies outside of the image of φ_μ , as discussed in

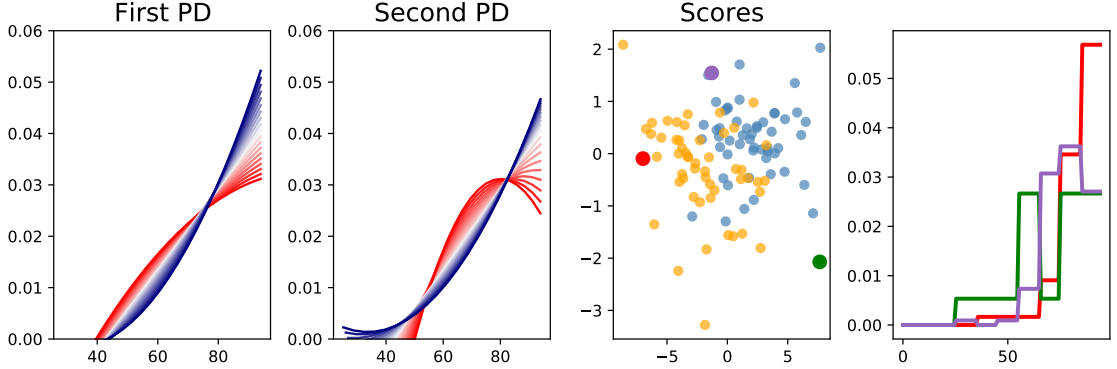


Figure 7: The first two panels show the variability along the first two principal directions (first and second panel), using the same visualization technique as in Figure 4. The third panel reports the scores of the projections on the two dimensional principal component (orange for women and blue for men) and the fourth panel shows three particular distributions, also highlighted in the third panel. In particular, the red distribution is the one of women in Vermont, the green one are males in Alaska and the purple one are women in West Virginia.

Section 3.4. As expected from the comparison between the metric projection and the pushforward operator in Figure 2, the fit to the data of the projected and log PCAs will be different. In particular, in this case we observe that the log PCA does a worse job in term of NRE , as shown in Figure 6(b), especially when the dimension increases. This behavior can be also in part explained by the complexity of the numerical routines needed to approximate the pushforward operator (required by the log PCA) where it is natural to expect some numerical errors.

More in general, as discussed also in Cazelles et al. (2018), we can conclude that the log PCA is not suited to study this particular data set because the L_2 PCA is different from the nested geodesic PCA (as testified by the GV_2 score). In fact, apart from the visual inspection of the L_2 principal directions – which are not guaranteed to span the log-principal components – not much can be obtained from the log PCA in this case, since it does not provide a consistent way of projecting data points on the principal component as pointed out in Section 3.4.

8. Numerical Illustrations for the Distribution on Distribution Regression

In this section, we propose a comparison between the Wasserstein projected and simplicial (see Appendix B) approaches when the task at hand is distribution on distribution regression, and show an application of the Wasserstein projected regression framework to a problem of wind speed forecasting.

8.1 Simulation Study

We consider two data generating processes as follows. In the first setting, data are generating from the Wasserstein regression: independent variables z_1, \dots, z_n are generated by considering quantile functions $F_{z_1}^-, \dots, F_{z_n}^-$ such that $F_{z_i} = \sum_{h=1}^{30} a_{ih}^{(z)} \psi_j^{(3)}$ where $\psi_1^{(3)}, \dots, \psi_{30}^{(3)}$ is a cubic spline basis over equispaced knots in $[0, 1]$ and $a_{i1}^{(z)} = 0$, $a_{i2}^{(z)} = \delta_{i1}$,

	First scenario	Second scenario
Wasserstein	$(4 \times 10^{-7}, 7 \times 10^{-8})$	$(5 \times 10^{-3}, 6 \times 10^{-3})$
Simplicial	$(0.9, 2.66)$	$(4 \times 10^{-4}, 5 \times 10^{-4})$

Table 1: Cross validation (leave one out) errors and standard deviations for the Wasserstein and Simplicial regression under the two simulated examples

$a_{ij}^{(z)} = a_{ij-1}^{(z)} + \delta_{ij-1}$, and $(\delta_{i2}, \dots, \delta_{i30}) \sim \text{Dirichlet}(1, \dots, 1)$. This data generating procedure ensures the $F_{zi}^-(0) = 0$, $F_{zi}^-(1) = 1$ and F_{zi}^- is monotonically increasing, cf. Proposition 4. The dependent variables $F_{y1}^-, \dots, F_{ym}^-$ are generated using the same spline expansion of the dependent variables and letting $\mathbf{a}_i^{(y)} = B\mathbf{a}_i^{(z)}$. B is a randomly generated matrix with rows $\mathbf{b}_1, \dots, \mathbf{b}_{30}$, and each \mathbf{b}_i is generated as follows: $b_{i1} \sim \text{Uniform}(0, 0.5)$ $b_{ij} = b_{ij-1} + \tilde{b}_{ij}$ and $\tilde{b}_{ij} \sim \text{Uniform}(0, 0.5)$, so that the coefficients $a_{ij}^{(y)}$ are monotonically non decreasing for each i and thus the F_{yi}^- 's can be considered quantile functions.

We compute the pushforward of the uniform distribution via numerical inversion and differentiation and obtain the pdf associated to each quantile function. Observe that this task is easier than approximating the pushforward of a generic μ through a generic f (as Cazelles et al. (2018) do) since the quantile functions are monotonic and we have simple expressions for all the quantities related to μ . Since the simplicial regression takes as input (a transformation of) the pdfs while the Wasserstein regression works directly on the quantile functions, and also due to the fact that numerical errors can be introduced in the data set during the inversion and differentiation, we consider as ground truth the pdfs and, for the Wasserstein approach, re-compute numerically the quantile functions.

In the second setting instead, we generate data from the simplicial regression model: independent variables z_1, \dots, z_n are generated by applying the inverse of the centered log ratio to a random spline expansion as follows. For each $i = 1, \dots, n$ let $\tilde{p}_{zi} = \sum_{j=1}^{30} a_{ij}^{(z)} \psi_j^{(3)}$ where the $\psi_j^{(3)}$'s are the same B-spline basis as in the previous setting. Here, the $a_{ij}^{(z)}$'s are generated iid from a Gaussian distribution with mean 0 and standard deviation 0.2. The dependent variables are generated by letting $\tilde{p}_{yi} = \sum_{j=1}^{30} a_{ij}^{(y)} \psi_j^{(3)}$ and $\mathbf{a}_i^{(y)} = B\mathbf{a}_i^{(z)}$, where B is a randomly generated 30×30 matrix with entries drawn iid from a standard normal distribution. Finally the pdfs p_{zi} (p_{yi}) are recovered by applying the inverse of the centered log ratio to \tilde{p}_{zi} (\tilde{p}_{yi}), see Appendix B for more details.

Note that under the second data generating process, both the dependent and independent distributions have support in $[0, 1]$ by construction, whereas under the first data generating process the independent variables might have a larger support. Thus, to fit the simplicial regression in the first scenario, as common practice (cf. Appendix B), we extend the support of all the distributions (both dependent and independent) to the smallest interval of the real line containing all the supports. This is done by adding a small term to the pdfs (in our example, 10^{-12}) and then renormalizing them.

For both examples, we simulated 100 observations and compared the projected-Wasserstein and simplicial regression using leave-one-out cross-validation. In particular, for both approaches we use $J = 20$ quadratic spline basis and choose the penalty term ρ in (16) through grid search. Table 1 shows the pairs of mean squared error and standard deviation of the cross validation, the metric to compare the ground truth and the prediction is the 2-Wasserstein distance. As one might expect, the Wasserstein regression performs better in the first scenario while the simplicial regression performs better in the second scenario. However, it is surprising how the Wasserstein geometry can capture (in terms of Wasserstein metric) dependence generated by a linear structure which we have shown to

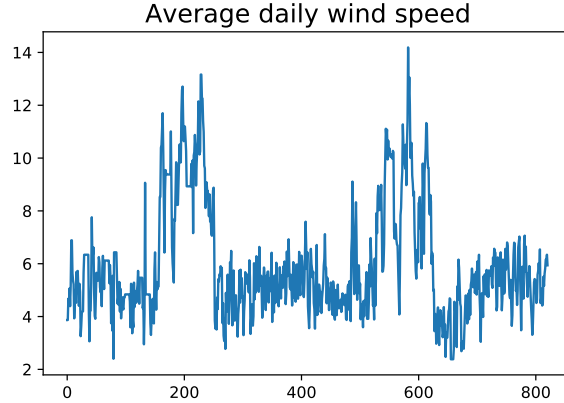


Figure 8: Daily average wind speed

be very different from the Wasserstein one, making the projected regression a promising tool for such inferential problems

8.2 Wind speed distribution forecasting from a set of experts

We consider the problem of forecasting the distribution of the wind speed nearby a wind farm from a set of experts. The data set is publicly available at www.kaggle.com/theforcecoder/wind-power-forecasting. In particular, data consists of measurements of the wind speed collected every ten minutes for a period of 821 days starting from the 31st December 2017. The daily average wind speed is shown in Figure 8.

We assume to have access to a set of *experts*, that is a set of trained models, that provide a probabilistic one-day-ahead forecast for the average wind speed. Here, our goal is to combine this set of experts and provide a point estimate of the wind speed distribution for the whole day, which can be helpful when planning the maintenance of the wind mills for instance.

Formally, let K denote the number of experts considered, F_{zij}^- is the quantile function associated to the probabilistic forecast of the average wind speed for day i given by expert $j = 1, \dots, K$; F_{yi}^- is the empirical quantile function of the wind speed for day i . In particular, we consider $K = 4$ experts built from the *Prophet* model by Facebook (Taylor and Letham, 2018) as follows: model $M1$ is the classical Prophet, without additional covariates or seasonality trends; model $M2$ includes the ambient temperature as covariate but not seasonality; model $M3$ includes a yearly seasonality and no covariates and model $M4$ includes both yearly seasonality and ambient temperature as covariate. The models are estimated using variational inference on rolling samples of 365 days and produce one day ahead probabilistic forecasts for the average wind speed. The final sample size corresponds to $n = 456$.

We consider a trivial extension of the distribution on distribution regression model in Section 5.2 as follows:

$$\mathbb{E}[F_{yi}^- | F_{zi1}^-, \dots, F_{ziK}^-] = \Pi_{L_2([0,1])^\uparrow} \left(\alpha + \sum_{j=1}^K \int_0^1 \beta_j(t, s) F_{zij}^-(t) dt \right) \quad (23)$$

	<i>R1</i>	<i>R2</i>	<i>R3</i>	<i>R2</i>	<i>RF</i>
MSE	(1.22 ± 1.32)	(1.19 ± 1.26)	(1.15 ± 1.07)	(1.24 ± 1.23)	(0.86 ± 0.82)

Table 2: Mean square prediction error \pm one standard deviation on the held-out test set.

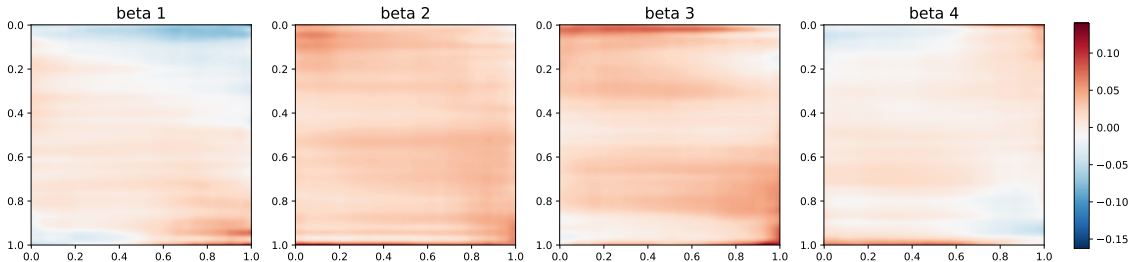


Figure 9: Estimates of the $\beta_i(t, s)$'s evaluated on $[0, 1]^2$. The variable t runs across columns, and variable s across rows

Having approximated all the functions through a B-spline expansion, the model reads

$$\mathbb{E}[\mathbf{a}_i^{(y)} \mid \mathbf{a}_{i1}^{(z)}, \dots, \mathbf{a}_{iJ}^{(z)}] = \Pi_{\mathbb{R}^{J \uparrow}} \left(\boldsymbol{\theta}_\alpha + \sum_{j=1}^K \Theta_{\beta_j} E \mathbf{a}_{ij}^{(z)} \right).$$

The procedure for estimating $\boldsymbol{\theta}_\alpha$ and $\Theta_{\beta_1}, \dots, \Theta_{\beta_K}$ is analogous to the one outlined in Section 5.2.

We compare the prediction performance of five distribution on distribution regression models. Models *R1* to *R4* are obtained by fitting model (23) using only one of the four experts, *M1* to *M4*, while the fifth model (*RF*) is the “full” model in (23) considering all the four experts. For this comparison, we perform a train-test split of the 456 days for which the experts produced the prediction, considering the last 100 days as test. We select hyperparameters (namely, the penalty coefficient ρ in (16) and whether to include or not the intercept term α) by a grid search cross validation on the training set, and compare the mean square error on the held-out test set. Results of the comparison are reported in Table 2. As expected, the model with the four predictors (*RF*) is the best performer. Interestingly, all the other models *R1-R4* perform similarly and present a much higher mean square error when compared to *RF*, thus suggesting that the best performance is achieved by combining the different experts together and no expert alone can be a good predictor. This is possibly explained by some experts being able to better forecast one scenario (for instance, light winds) and other experts being able to better forecast other scenarios.

We conclude with some descriptive analysis. Figure 9 shows the point estimates for the coefficients β_j . We can interpret as highly influential for the regression the areas of the β_j 's with high absolute value, and as negligible area with values close to zero.

We can highlight some differences among the coefficients in Figure 9. In particular, model *M1*, seems influent when predicting the tails of the distribution, in particular with negative weights for the left tail and positive weights for the right tail. Model *M2*, seems to be affecting all the steps of the prediction and in particular to be model affecting the most the median of the distribution. Model *M3*, appears to be, with *M2*, the most important model for the prediction: the absolute value in the corresponding regressor β_3 is often very high and with noticeable peaks corresponding to areas predicting the left tail and

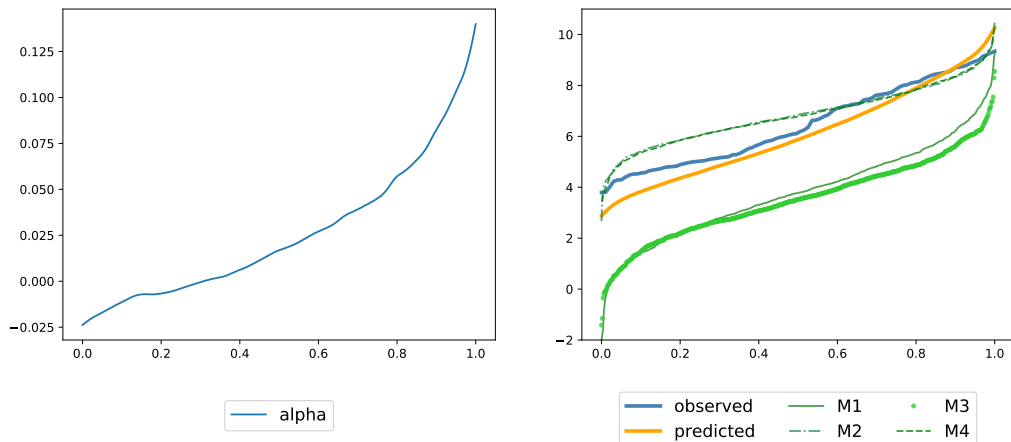


Figure 10: Estimate of α (left) and prediction of one F_y^- of the test set (right). In the right panel, the blue line corresponds to the empirical quantile function, the orange one to the prediction from RF and the green ones to the average wind predictions obtained from the experts $M1$ - $M4$.

towards the right tail. Finally, the regressor corresponding to $M4$ has very low values thus resulting in minor importance in terms of regression influence.

Interestingly, the experts providing the most precious inputs to our regression model are $M2$ and $M3$, that incorporate only the seasonality effect and the temperature covariate respectively, while $M4$, which incorporates both, seems to be less important. Hence, the regression model in (23) finds more effective combining experts trained on different covariates than correcting an expert already trained on all the covariates. In particular, our insight is that $M2$ is responsible for centering the median of the output distribution. The tails of the distribution seem to need also the contribution of seasonality data, given by $M3$. Finally, we also observe that the left tail of the wind distribution seems the most difficult to be predicted, needing very high positive and negative weights across different models, to be obtained.

9. Discussion

In this paper, we propose a novel class of *projected* statistical methods for distributional data on the real line, focusing in particular on the definition of a *projected* PCA and a *projected* linear regression. By investigating the weak Riemannian structure of the Wasserstein space and the transport maps between probability measures, we represent the Wasserstein space as a closed convex cone inside an Hilbert space.

Similar to *log* methods, our models exploit the possibility to map data into a linear space to perform statistics in an *extrinsic* fashion. However, instead of using operators like the *exp* map or a some kind of boundary projection to return to the Wasserstein space, we rely on a metric projection operator that is more respectful of the underlying metric.

By choosing as base point the uniform measure on $[0, 1]$, we are able to efficiently approximate the metric projection operator so that our models combine the ease of implementation of *extrinsic* methods while retaining a performance similar to the one of *intrinsic* methods. Further, through a quadratic B-spline approximation, we can greatly reduce the dimensionality of the optimization problems involved, resulting in fast empir-

ical methods. As a byproduct of this approach, we also derive faster numerical routines for the *geodesic* PCA in Bigot et al. (2017).

We study asymptotic properties of the proposed methods, concluding that, under reasonable regularity assumptions, our *projected* models provide consistent estimates and that the B-spline approximation error becomes negligible. We showcase our approach in several simulation studies and using two real world data sets, comparing our models to *intrinsic* and *extrinsic* ones and to the *simplicial* approach in Hron et al. (2014), concluding that the *projected* PCA and regression constitute a valid candidate for performing inference on a data set of distributions.

Although our *projected* framework was proven to be viable in many practical situations, some care must be taken when adopting it, especially when performing PCA. In fact, the *extrinsic* nature of our method might not fit every data set, in which case a more computationally demanding *intrinsic* PCA might be preferred, see for instance Appendix D.1 for an example where the *projected* principal directions are not interpretable. On top of that, performing PCA in the Wasserstein space requires more attention than performing the usual Euclidean PCA: as pointed out in Appendix D.2, since principal components are not linear subspaces, decomposing the variance along the directions (i.e., looking at the scores) must be done carefully, and making sure that the directions are indeed interpretable. To assist practitioners, in Section 7.2 we have also proposed two scores that quantify the interpretability of the principal directions and the discrepancy between the *nested* and *projected* principal components.

Several extensions and modifications of our approach are possible. One possibility is to extend our framework to encompass more models, such as generalized linear models and independent component analysis. Although this should be straightforward in theory, the numerical computations could become more burdensome. Furthermore, as an alternative to our approach based on B-splines approximation, one could use such B-spline expansion only to approximate the metric projection operator. Another interesting line of research would consist in building hybrid approaches (as anticipated in Section 7.2) to analyze distributions in the Wasserstein space, using both *extrinsic* and *intrinsic* methods to exploit the advantages of both worlds, while mitigating the disadvantages. We also think that a deeper comparison between the Wasserstein and the simplicial geometries could help practitioners in choosing between them.

Finally, as pointed out by an anonymous referee, extensions to encompass measures supported on \mathbb{R}^d , $d > 1$, are of great interest. This is surely a very challenging problem, due to the geometric structure of $\mathcal{W}_2(\mathbb{R}^d)$. We identify three main obstacles in this sense. First, the map onto the tangent space is not an isometry because the Wasserstein space is curved. Second, we lose the nice characterization of the tangent space and of the image of \log_μ , so that the metric projection operator becomes harder to derive. Third, the computational cost would greatly increase due to the need of numerically approximating the transport maps needed to compute the distances.

Acknowledgments

We thank two anonymous referees for their detailed and helpful reviews, which helped us improving the quality and the clarity of this work. We also thank Riccardo Scimone for helpful feedback and comments on an earlier draft of this paper and Federico Bassetti, Alessandra Guglielmi and Piercesare Secchi for helpful discussions.

Appendix A. Proofs

Assumptions on x_0 .

Let $B_\varepsilon(x_0) = \{x \in H \mid \|x - x_0\| < \varepsilon\}$, a ball of radius ε in H . Given a set C , we refer to $\text{aff}(C)$ as the smallest affine subset containing C , found as the intersection of all affine subspaces containing C . Similarly $\mathcal{H}(C)$ is the convex hull of C , the smallest convex subset of H containing it. The relative interior of a set C is defined as its interior considering as ambient space $\text{aff}(C)$: $\text{relint}(C) = \{x \in C \mid \exists B_\varepsilon(x_0) \text{ such that } B_\varepsilon(x_0) \cap \text{aff}(C) \subset C\}$.

Throughout our paper we assume that the random variable \mathcal{X} is such that (i) there exists $x_0 = \mathbb{E}[\mathcal{X}]$ and (ii) $x_0 \in \text{relint}(\mathcal{H}(\text{supp}(\mathcal{X})))$ where $\text{supp}(\mathcal{X})$ is the support of \mathcal{X} . These assumptions are indeed quite natural and require that the distribution of \mathcal{X} has a well defined barycenter, which is not in a “degenerate” position with respect to the convex hull of its support, which may happen in infinite dimensional Hilbert Spaces. See, for instance, Berezin and Miftakhov (2019) for an example of distributions not verifying this second assumption.

Proof of Lemma 1.

The proof is divided in two steps. First, we prove that $(x_0 + \text{Sp}(U_k)) \cap X$ has dimension k . Then, we show that $U_X^{x_0, k} = (x_0 + \text{Sp}(U_k)) \cap X$. Without loss of generality, for ease of notation, we perform an affine change of variable so that $x_0 = 0$, but, with a slight abuse of notation, we keep denoting with \mathcal{X} and X the transformed random variable and the convex cone respectively.

To prove the first part, let $\mathcal{H}(\mathcal{X})$ be the convex hull of the support of \mathcal{X} and $\text{aff}(\mathcal{H}(\mathcal{X})) = K$ be the smallest affine subset of H containing $\mathcal{H}(\mathcal{X})$. We know by assumption that there is an open ball in K which contains $x_0 = 0$ and is contained in $\mathcal{H}(\mathcal{X})$. Moreover, for every $k \leq \dim(K)$, $\text{Sp}(U_k) \subset K$. Note that we can clearly suppose $k \leq \dim(K)$, otherwise principal components analysis is useless. With this assumption, since $x_0 = 0$ is in the relative intern of $\mathcal{H}(\mathcal{X})$, we have $k = \dim(\text{Sp}(U_k) \cap \mathcal{H}(\mathcal{X})) \leq \dim(\text{Sp}(U_k) \cap X) \leq k$.

Now we prove that a $(k, 0)$ -projected principal component is given by $\text{Sp}(U_k) \cap X$. To prove this, let C^* be a $(k, 0)$ -projected principal component and $A^* = A \cap X$, with $A = \text{Sp}(U_k)$: we know (i) $x_0 = 0 \in A^*$, (ii) $\dim(A^*) = k$ by definition and (iii) $A^* \subseteq \Pi_X(A)$, so we have $A^* \subset C^*$.

Since $\dim(C^*) = k$ there is C linear subspace of dimension k such that $C^* \subset C$. Consider $C' = C \cap X$: clearly $C^* \subset C'$, so that $A^* \subset C^* \subset C'$. Moreover, $A^* \subset C'$, which implies $A \cap X \subset C \cap X$ and thus $\text{Sp}(A \cap X) \subset \text{Sp}(C \cap X)$. The proof is concluded if $\dim(\text{Sp}(A \cap X)) = \dim(\text{Sp}(C \cap X)) = k$. In fact, in this case $A = \text{Sp}(A \cap X)$ and $C = \text{Sp}(C \cap X)$ which means that $A \subset C$ and since $\dim(A) = \dim(C) = k$, A and C coincide, proving $A^* = C^*$.

To prove this final claim, observe that $\dim(\text{Sp}(A \cap X)) < k$ implies $\dim(A \cap X) < k$, which contradicts the proof of the first part of this Lemma. Similarly, $\dim(\text{Sp}(C \cap X)) = k$ since $\dim(C^*) = k$ by hypothesis. ■

Proof of Proposition 1.

The fact that $\|\Pi_{U_X^{x_0, k}}(x) - x\| \geq \|\Pi_{U_X^{x_0, k+1}}(x) - x\|$ follows easily by noticing that $U_X^{x_0, k} \subset U_X^{x_0, k+1}$.

Now, to prove that $\|\Pi_{U_X^{x_0, k}}(x) - x\| \rightarrow 0$ as k increases, we first notice that, by the properties of the principal components in H we have $\Pi_{\text{Sp}(U_k)}(x - x_0) \xrightarrow{k} x - x_0$ for every

$x \in X$, which implies $\|\Pi_{Sp(U_k)+x_0}(x) - x\| \rightarrow 0$. Then, denote $x_1 = \Pi_{U_X^{x_0,1}}(x)$ and let r_k be the line between x_1 and x . Let:

$$x_k = \arg \min_{x' \in r_k \cap Sp(U_k)+x_0} \|x' - x\|$$

We clearly have $x_k \rightarrow x$. Finally, by convexity we know $x_k \in U_X^{x_0,k}$, which implies $\|\Pi_{U_X^{x_0,k}}(x) - x\| \leq \|x_k - x\| \rightarrow 0$. ■

Proof of Proposition 2.

Again, without loss of generality, for ease of notation, we perform an affine change of variable so that $x_0 = 0$, but, with a slight abuse of notation, we keep denoting with \mathcal{X} and X the transformed random variable and convex cone respectively.

We start by noticing that being Π_k the orthogonal projection onto a subspace, $x - \Pi_k(x) \perp Span(U_k)$ and thus for $v \in Span(U_k)$:

$$\|x^* - v\|^2 = \|x^* - \Pi_k(x^*)\|^2 + \|\Pi_k(x^*) - v\|^2$$

Then

$$\arg \min_{v \in U_X^{0,k}} \|x^* - v\| = \arg \min_{v \in Sp(U_k) \cap X} \|\Pi_k(x^*) - v\|$$

and the result follows. ■

Proof of Proposition 4.

1. As shown in the supplementary of Pya and Wood (2015) by standard B-spline formulas we obtain that given $f(x) = \sum_{j=1}^J a_j \psi_j^k(x)$, then $f'(x) = \sum_{j=1}^J (a_j - a_{j-1}) \cdot \psi_j^{k-1}(x)$. Being the B-spline basis function nonnegative by definition, we obtain the result.
2. With $k = 2$, $f'(x)$ on the interval $[x_{j+1}, x_j]$ has the following expression:

$$\frac{x - x_j}{x_{j+1} - x_j} \cdot (\alpha_j - \alpha_{j-1}) + \frac{x_{j+1} - x}{x_{j+1} - x_j} \cdot (\alpha_{j-1} - \alpha_{j-2})$$

so:

$$\lim_{x \rightarrow x_{j+1}^-} f'(x) = \alpha_j - \alpha_{j-1}$$

and the result follows. ■

Proof of Proposition 5 and 6.

We report here Propositions 3.3 and 3.4 of Bigot et al. (2017), with the notation adapted to our manuscript. In the following H is a separable Hilbert space, X is a closed convex subset of H , \mathcal{X} is an X -valued square integrable random variable, x_0 a point in X and $k \geq 1$ an integer.

Proposition 10 *Let $U^* = \{u_1^*, \dots, u_k^*\}$ be a minimizer over orthonormal sets U of H of cardinality k , of $D_X^{x_0}(\mathcal{X}, U) := \mathbb{E}d^2(\mathcal{X}, (x_0 + Sp(U)) \cap X)$, then $U_X^{x_0} := (x_0 + Sp(U)) \cap X$ is a (k, x_0) -global principal component of \mathcal{X} .*

Proposition 11 *Let $U^* = \{u_1^*, \dots, u_k^*\}$ be an orthonormal set such that $U_i^* = \{u_1^*, \dots, u_i^*\}$ is a minimizer of $D_X^{x_0}(\mathcal{X}, U)$ over the orthonormal sets of cardinality “ i ” such that $U \supset U_{i-1}^*$; then $U_X^{*x_0}$ is a (k, x_0) -nested principal convex component of \mathcal{X} .*

Applying Propositions 10 and 11 we can obtain equivalent definitions of geodesic and nested PCA as optimization problems in $L_2([0, 1])$. If we fix $J \in \mathbb{N} > 0$ and a quadratic B-spline basis $\{\psi_j\}_{j=1}^J$, we can use Propositions 10 and 11 with $X = L_2([0, 1])^{J\uparrow}$ and $H = L_2([0, 1])^J$. Thanks to Remark 7 we obtain the results. ■

Proof of Proposition 7.

Let $S_J = \sum_{j=1}^J \lambda_j^{(J)} \psi_j^{(J)}$ and its derivative $s_J = \sum_j (\lambda_j^{(J)} - \lambda_{j-1}^{(J)}) \tilde{\psi}_j^{(J)}$ where $\tilde{\psi}_j^{(J)}$ denotes the linear spline basis on the same equispaced grid in $[0, 1]$.

Let $f_\mu^- = (F_\mu^-)'$, of course it can be seen that f_μ^- is non-negative. Moreover, it is obvious that $f_\mu^- \in W_2^\infty([0, 1])$. Then, from De Boor and Daniel (1974) we get that there exist s_J such that $\|s_J - f_\mu^-\|_\infty \leq C \|D^2 f_\mu^-\|_\infty J^{-2}$, where C is a constant depending on the interval $[0, 1]$ but not on n .

Hence, we can determine the coefficients $\{\lambda_j^{(J)}\}$, starting from the spline s_J , up to a translation factor.

We fix a particular set of coefficients by letting $S_J(0) = \lambda_1^{(J)} = F_\mu^-(0)$ for each J . So that:

$$S_J(x) - F_\mu^-(x) = \int_0^x s_J(t) dt - \int_0^x f_\mu^-(t) dt - S_J(0) + F_\mu^-(0) = \int_0^x s_J(t) - f_\mu^-(t) dt$$

By using the previous result, the integral we have that $S_J(x) - F_\mu^-(x) \leq C J^{-2}$ for all x which proves the proposition. ■

Proof of Proposition 8.

By the Assumptions in Section 6.2.1 and Remark 10 there exists a ball B_K in $W_3^\infty([0, 1])$ of radius K for some $K > 0$, such that each F_i^- can be ε -approximated by $\tilde{F}_i^- \in W_3^\infty([0, 1])$ with $\tilde{F}_i^- \in B_K$. We can suppose that also the eigenvectors of the covariance operator of the generating process belong to such sphere, otherwise we just increase its radius of some finite amount.

By Proposition 7 we can choose a spline basis (that is, a number of elements $J > 0$), such that we get a ε -uniformly good approximation of B_K (and thus we can 2ε -approximate its L_2 closure). To lighten notation, thanks to Remark 7 we deliberately confuse $\mathbb{R}^{J\uparrow}$ and the space monotone B -splines with J basis functions, the inner product we are referring to will always be clear by looking at its entries.

Now consider the following inequalities, with \mathbf{a}_i^J obtained as 2ε approximations of F_i^- , $\mathbf{w}^J \in \mathbb{R}^J$, $w \in L_2([0, 1])$:

$$\left| \frac{1}{n} \sum_i \langle F_i^-, w \rangle^2 - \frac{1}{n} \sum_i \langle \mathbf{a}_i^J, \mathbf{w}^J \rangle^2 \right| \leq \frac{1}{n} \left| \sum_i \langle F_i^-, w \rangle^2 - \sum_i \langle \mathbf{a}_i^J, w \rangle^2 + \sum_i \langle \mathbf{a}_i^J, w \rangle^2 - \sum_i \langle \mathbf{a}_i^J, \mathbf{w}^J \rangle^2 \right|,$$

where the inner product $\langle \mathbf{a}_i^J, w \rangle$ is to be intended as the L_2 inner product between the spline function with coefficients \mathbf{a}_i^J and the L_2 function w . Consider now:

$$\begin{aligned} \frac{1}{n} \sum_i (\langle F_i^-, w \rangle^2 - \langle \mathbf{a}_i^J, w \rangle^2) &= \\ \frac{1}{n} \sum_i (\langle F_i^-, w \rangle - \langle \mathbf{a}_i^J, w \rangle)(\langle F_i^-, w \rangle + \langle \mathbf{a}_i^J, w \rangle) &= \\ \frac{1}{n} \sum_i \langle F_i^- - \mathbf{a}_i^J, w \rangle \langle F_i^- + \mathbf{a}_i^J, w \rangle &\leq \\ \frac{1}{n} \sum_i \left| \langle F_i^- - \mathbf{a}_i^J, w \rangle \right| \cdot \left| \langle F_i^- + \mathbf{a}_i^J, w \rangle \right| &\leq \\ \frac{1}{n} \sum_i 2\varepsilon \|w\|^2 2K &= 4\varepsilon K \|w\|^2 \end{aligned}$$

Similarly:

$$\left| \frac{1}{n} \sum_i (\langle \mathbf{a}_i^J, w \rangle^2 - \langle \mathbf{a}_i^J, \mathbf{w}^J \rangle^2) \right| \leq \|\mathbf{a}_i^J\|^2 \cdot \|w - \mathbf{w}^J\| \cdot (\|w\| + \|\mathbf{w}^J\|)$$

We know that a solution to the problem $\max_{\|w\|_{L_2}=1} \frac{1}{n} \sum_i \langle F_i^-, w \rangle^2$ is given by the first eigenfunction \hat{w} of the covariance operator of the empirical process. Now we are in the condition to apply results in Dauxois et al. (1982), or in Qi and Zhao (2011) (with $\alpha \rightarrow 0$) to conclude that \hat{w} converges to the first eigenfunction \bar{w} of the covariance operator of the process that generates F_i^- . By hypothesis, such eigenfunction \bar{w} lies in B_K and thus can be approximated with our fixed spline basis. Thus for high enough n , also \hat{w} can be approximated up to 2ε .

Let $\mathbf{a}_{\hat{w}}$ be the coefficients of the spline expansion of \hat{w} spline approximation, that is, $\|w - \mathbf{a}_w\| \leq 2\varepsilon$. Observe that $\left| \|\hat{w}\|_2 - \|\mathbf{a}_{\hat{w}}\|_E \right| \leq 2\varepsilon$, just as $\|\mathbf{a}_i^J\| \leq K + 2\varepsilon$. Thus, up to adding another ε to the approximation error $\|\hat{w} - \mathbf{a}_{\hat{w}}\|$, we can suppose $\|\mathbf{a}_{\hat{w}}\|_2 = 1$. Hence:

$$\left| \frac{1}{n} \sum_i (\langle \mathbf{a}_i^J, \hat{w} \rangle^2 - \langle \mathbf{a}_i^J, \mathbf{a}_{\hat{w}} \rangle^2) \right| \leq (K + 2\varepsilon) \cdot 3\varepsilon \cdot 2$$

Which leads to:

$$\left| \max_{\|w\|_{L_2}=1} \sum_i \langle \mathbf{a}_i^J, w \rangle^2 - \max_{\|\mathbf{w}^J\|_E=1} \sum_i \langle \mathbf{a}_i^J, \mathbf{w}^J \rangle^2 \right| \leq (K + 2\varepsilon) \cdot 3\varepsilon \cdot 2$$

Finally, combining the above results and the fact that $|\max f - \max g| \leq \max |f - g|$ for any pair of real valued functions f and g , we obtain:

$$\begin{aligned} \left| \max_{\|w\|_{L_2}=1} \frac{1}{n} \sum_i \langle f_i, w \rangle^2 - \max_{\|\mathbf{w}^J\|_E=1} \frac{1}{n} \sum_i \langle \mathbf{a}_i^J, \mathbf{w}^J \rangle^2 \right| &\leq \\ \max_{\|w\|_{L_2}=1} 4\varepsilon K \|w\| + (K + 2\varepsilon) \cdot 6\varepsilon &\leq 6\varepsilon K (1 + 2\varepsilon) \end{aligned}$$

Thus for instance if we ask that $\varepsilon < 1$, we obtain the desired result with $D = 18 \cdot K$. Consistency follows since $\|\mathbf{a}_{\hat{w}} - \bar{w}\| \leq \|\mathbf{a}_{\hat{w}} - \hat{w}\| + \|\hat{w} - \bar{w}\|$. \blacksquare

Proof of Lemma 2.

Since for any $x \in X$ we have $\Pi_{\mathbb{R}^{J \uparrow}}(x) \rightarrow x$, for any $v \in H$:

$$\|v - \Pi_{\mathbb{R}^{J \uparrow}}(v)\| \leq \|v - \Pi_{\mathbb{R}^{J \uparrow}}(\Pi_X(v))\| \leq \|v - \Pi_X(v)\| + \|\Pi_X(v) - \Pi_{\mathbb{R}^{J \uparrow}}(\Pi_X(v))\|$$

which implies $\Pi_{\mathbb{R}^{J \uparrow}}(v) \rightarrow \Pi_X(v)$. Consider now $\beta_n \rightarrow \beta$ in H ; we have the inequality:

$$\|\Pi_{\mathbb{R}^{J \uparrow}}(\beta_n) - \Pi(\beta)\| \leq \|\Pi_{\mathbb{R}^{J \uparrow}}(\beta_n) - \Pi_X(\beta_n)\| + \|\Pi_X(\beta_n) - \Pi_X(\beta)\|$$

the first term of the right hand side of the inequality can be sent to 0 by increasing J , the other by increasing n . \blacksquare

Proof of Proposition 9.

We call a_i the spline coefficients associated to x_i and b_i the ones associated to y_i . Again we deliberately confuse the spaces where the coefficients live to lighten the notation. Since the penalty term does not depend on the data, we have:

$$\begin{aligned} & \frac{1}{n} \left| \sum_i \|y_i - \langle x_i, B^T AB \rangle\|^2 - \sum_i \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|^2 \right| = \\ & \frac{1}{n} \left| \sum_i (\|y_i - \langle x_i, B^T AB \rangle\|^2 - \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|^2) \right| \leq \\ & \frac{1}{n} \sum_i \|y_i - \langle x_i, B^T AB \rangle\|^2 - \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|^2 \end{aligned}$$

Now, since

$$\begin{aligned} & \left| \|y_i - \langle x_i, B^T AB \rangle\|^2 - \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|^2 \right| = \\ & \left| (\|y_i - \langle x_i, B^T AB \rangle\| - \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|) \times \right. \\ & \left. (\|y_i - \langle x_i, B^T AB \rangle\| + \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|) \right| \end{aligned}$$

Then for some constant K depending on the bounds in the Assumptions, we get:

$$\begin{aligned} & \left| \|y_i - \langle x_i, B^T AB \rangle\|^2 - \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|^2 \right| \leq \\ & \|y_i - \langle x_i, B^T AB \rangle - b_i + \langle a_i, B^T AB \rangle_{L_2([0,1])}\| 2K = \\ & (\|y_i - b_i\| + \langle a_i - x_i, B^T AB \rangle) 2K \end{aligned}$$

Thus, if J is such that we have ε -approximations of the data, by Cauchy-Schwartz we obtain:

$$\frac{1}{n} \left| \sum_i \|y_i - \langle x_i, B^T AB \rangle\|^2 - \sum_i \|b_i - \langle a_i, B^T AB \rangle_{L_2([0,1])}\|^2 \right| \leq K' \cdot \varepsilon$$

for some K' constant.

Thanks to the results in Prchal and Sarda (2007), for any $\varepsilon > 0$, if the number of samples is big, $\hat{\Theta}$ and $\hat{\Theta}_J$ exist with probability $1 - \varepsilon$ and are unique. Since the value of the minimization problem the solve are arbitrarily close, then the minimizers converge in $\mathbb{R}^{J \times J}$ with the metric given by the spline basis. \blacksquare

Strong convergence implies semi-norm convergence.

Let \mathcal{Z} be an H -valued random variable and $\mathcal{C}_{\mathcal{Z}}$ the covariance operator associated to \mathcal{Z} , that is:

$$(\mathcal{C}_{\mathcal{Z}}f)(s) = \int_{[0,1]} \text{cov}(\mathbf{x}(s), \mathbf{x}(t))f(t)dt.$$

In the following, we denote with $\|\cdot\|_{L_2}$ the $L_2([0,1]^2)$ norm. Further, recall that $\|\text{cov}(\mathcal{Z}(s), \mathcal{Z}(t))\|_{L_2([0,1]^2)} = \mathbb{E}\|\mathcal{Z}\|^2$. We want to look at the behavior of $\|\widehat{\beta}_{\text{PS}} - \widehat{\beta}_J\|_{\mathcal{C}_{\mathcal{Z}}}$.

$$\begin{aligned} \int_{[0,1]} \langle \mathcal{C}_{\mathcal{Z}}(\widehat{\beta}_{\text{PS}}(s,t) - \widehat{\beta}_J(s,t)), \widehat{\beta}_{\text{PS}}(s,t) - \widehat{\beta}_J(s,t) \rangle dt &\leq \\ \|\mathcal{C}_{\mathcal{Z}}(\widehat{\beta}_{\text{PS}}(s,t) - \widehat{\beta}_J(s,t))\|_{L_2} \cdot \|\widehat{\beta}_{\text{PS}}(s,t) - \widehat{\beta}_J(s,t)\|_{L_2} &\leq \\ \mathbb{E}\|\mathbf{x}\|^2 \cdot \|\widehat{\beta}_{\text{PS}}(s,t) - \widehat{\beta}_J(s,t)\|_{L_2} \cdot \|\widehat{\beta}_{\text{PS}}(s,t) - \widehat{\beta}_J(s,t)\|_{L_2}. \end{aligned}$$

So $\|\widehat{\beta}_{\text{PS}} - \widehat{\beta}_J\|_{\mathcal{C}_{\mathcal{Z}}} \leq M \cdot \|\widehat{\beta}_{\text{PS}} - \widehat{\beta}_J\|_{L_2}^2$ for some constant M . Thus $\|\cdot\|_{L_2}$ convergence implies $\|\cdot\|_{\mathcal{C}_{\mathcal{Z}}}$ convergence.

Appendix B. The simplicial approach

The simplicial approach to distributional data analysis is based on the definition of Bayes space $\mathcal{B}^2(I)$ (Egozcue et al., 2006). Formally, let $I \subset \mathbb{R}$ a closed interval, the Bayes spaces $\mathcal{B}^2(I)$ is defined the equivalence class of probability densities $p(x)$ on I (that is $p(x) \geq 0$ and $\int_I p(x)dx = 1$) with square integrable logarithm.

The Bayes space is endowed with a linear space starting from the definition of the perturbation and powering operators, that are analogous to the sum and multiplication times a scalar, and inner product. Moreover Menafoglio et al. (2014) defines an isometric isomorphism between $\mathcal{B}^2(I)$ and $L_2([0,1])$ through the so-called centered log ratio (clr) map defined as

$$\tilde{p}(x) := \text{clr}(p)(x) = \log(p(x)) - \frac{1}{b-a} \int_a^b \log p(t)dt \quad (24)$$

for every $p \in \mathcal{B}^2(I)$. The inverse map is defined as

$$p(x) = \text{clr}^{-1}(\tilde{p})(x) = \frac{\exp(\tilde{p}(x))}{\int_I \exp(\tilde{p}(x))dx}$$

Thus, it is possible to define a *simplicial* PCA and *simplicial* regression on the Bayes space starting from the clr map. In particular, let p_1, \dots, p_n be observed densities on the interval I and let $\tilde{p}_i = \text{clr}(p_i)$. Denote with $\tilde{w}_1, \dots, \tilde{w}_k$ the first k principal directions estimated from the \tilde{p}_i 's, then a k dimensional simplicial principal component is the span of $\{w_i = \text{clr}^{-1}(\tilde{w}_i)\}_{i=1}^k$ in $\mathcal{B}^2(I)$.

Similarly, for pdfs $\{(p_z, p_y)_i\}_{i=1}^n$ a simplicial regression model is defined starting from the clr transformed variables. Let $\tilde{\Gamma}$ denote a functional regression model in L_2 for variables $\{(\tilde{p}_z, \tilde{p}_y)_i\}_{i=1}^n$, then the simplicial regression states:

$$\mathbb{E}[p_{yi} | p_{zi}] = \text{clr}^{-1} \left(\tilde{\Gamma}(\tilde{p}_{zi}) \right).$$

Apart from the different geometries of the Wasserstein and Bayes space, which are discussed in Sections 7 and 8, we can highlight one particular drawback from the simplicial approach, which we believe poses a significant limit to its usefulness. In fact, the main assumption is that all the pdfs p_i share the same support, which might not be the case

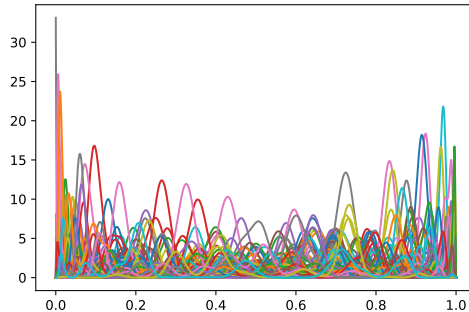


Figure 11: Example of data set from (26)

(for instance, it is not the case for our example in Section 8.2). In practice, one may circumvent this need by either “padding” all the pdfs to the same support, i.e considering

$$\bar{p}_i(x) \propto p_i(x) + \varepsilon \mathbb{I}[x \in I], \quad (25)$$

where $\mathbb{I}[\cdot]$ denotes the indicator function, and the proportionality is due to the need of re-normalizing the \bar{p}_i 's so that they integrate to 1. Another approach could consist in considering I as the intersection of all the supports of the different p_i 's let truncate all the pdfs to the shared interval I .

Both approaches present undesired side effects that can greatly alter the results. The second approach might end up with a very small interval I , so that a lot of information is lost due to this pre-processing step. The drawback of the first approach instead is due to numerical instability. In fact, one would like ε in (25) to be small in order not to corrupt the true signal, given by p_i . However, considering the transformation in (24) having a small ε would cause the \tilde{p}_i to present some extreme values (negative) in correspondence to ε . Performing PCA on a data set processed in this way would greatly alter the results, as most of the variability of the \tilde{p}_i 's would be masked by a difference in their support.

Appendix C. Additional Simulations

C.1 Sensitivity Analysis to the Number of Basis Functions

In this simulation, we show how the number of B-spline basis functions affects the inference in our projected PCA and in the simplicial one. In this Scenario, the probability measures are simulated as mixture of beta densities, also known as Bernstein polynomials, as follows:

$$p_i(x) = \sum_{j=1}^K w_{ij} \beta(x; j, K - j) \quad (26)$$

$$\mathbf{w}_i \sim \text{Dirichlet}_K(0.01)$$

Where $\beta(x; a, b)$ denotes the density of a beta distributed random variable with parameters (a, b) evaluated in x . By definition, the p_i s generated from (26) have a fixed support $I = [0, 1]$. See Figure 11.

In this setting instead, we let μ_i in (20) be the probability measure associated to p_i and not its smoothed version. Hence, in addition to the amount of information lost during the PCA another factor comes into play: the amount of information that is lost due to the B-spline representation.

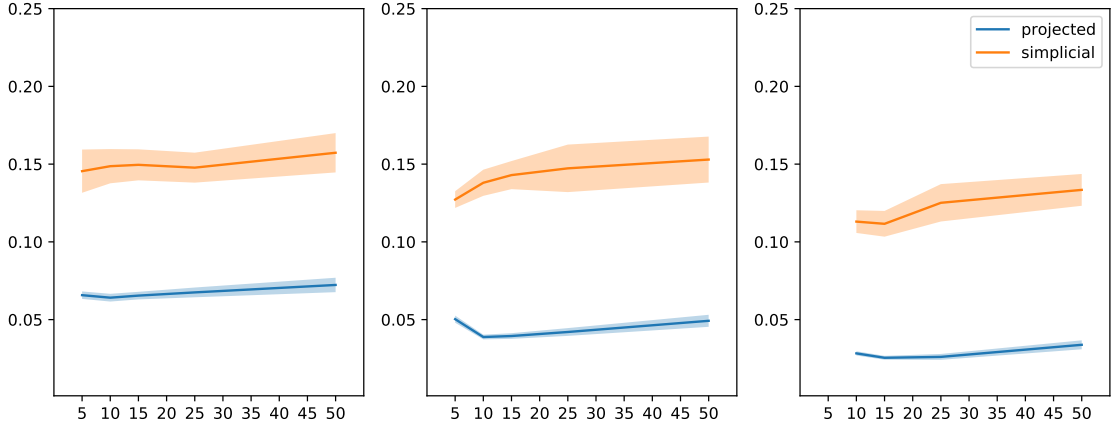


Figure 12: Results for the third scenario. All the panels show the reconstruction error as a function of the number of the spline basis functions. From left to right the results are obtained using the 2, 5 and 10 dimensional PCA. The solid lines represent the mean of 10 independent runs on independent data sets from (26) and the shaded area represent \pm one standard deviation.

Figure 12 shows the results. We can see that the reconstruction errors decrease when the dimension of the principal component increases both for the simplicial and projected PCA. Moreover, as the number of B-spline basis increase, the performance tend to get a little bit worse for both the approaches. We believe that this is due to an increased variance in the B-spline estimation of the quantile functions and (clr of) pdfs. In fact, computing the spline approximation for a single function amounts to solving a linear regression problem and increasing the dimension of the B-spline basis corresponds to increasing the number of regressors. Hence, letting B the matrix with columns ψ_1, \dots, ψ_J (evaluated on a grid), the variance of the OLS estimate of the coefficients \mathbf{a} is proportional to $(B^T B)^{-1}$. When increasing the number of B-splines, the entries in $B^T B$ become closer to zero, since the support of each of the spline basis becomes smaller. This leads to smaller precision (and higher variance) in the estimator for \mathbf{a} .

Another interesting thing to notice is that the simplicial PCA exhibits a much larger variance in the reconstruction error. This is possibly due to the different degree of smoothness of the quantile functions and of the pdfs. As the quantile functions are smoother than the pdfs, their B-spline basis expansion should have lower variance and be more similar to the true quantiles.

C.2 Empirical Verification of Consistency Results and Choosing J

In this section, we provide additional simulations to verify the consistency results established in Section 6.

For the PCA, we consider the two data generating processes in equations (19) (Gaussian) and (21) (DPM). First, first we fix $J = 20$ spline basis (as we do throughout Section 7) and let n increase. Then, we also let J increase linearly with n . We estimate the “true” principal directions by simulating 10^5 observations and using 2500 elements in the B-spline basis. Then, for any choice of n and J we generate another data set and compute the corresponding first two principal directions via the projected PCA and compute the L_2 norm between the “true” directions and the estimated ones.

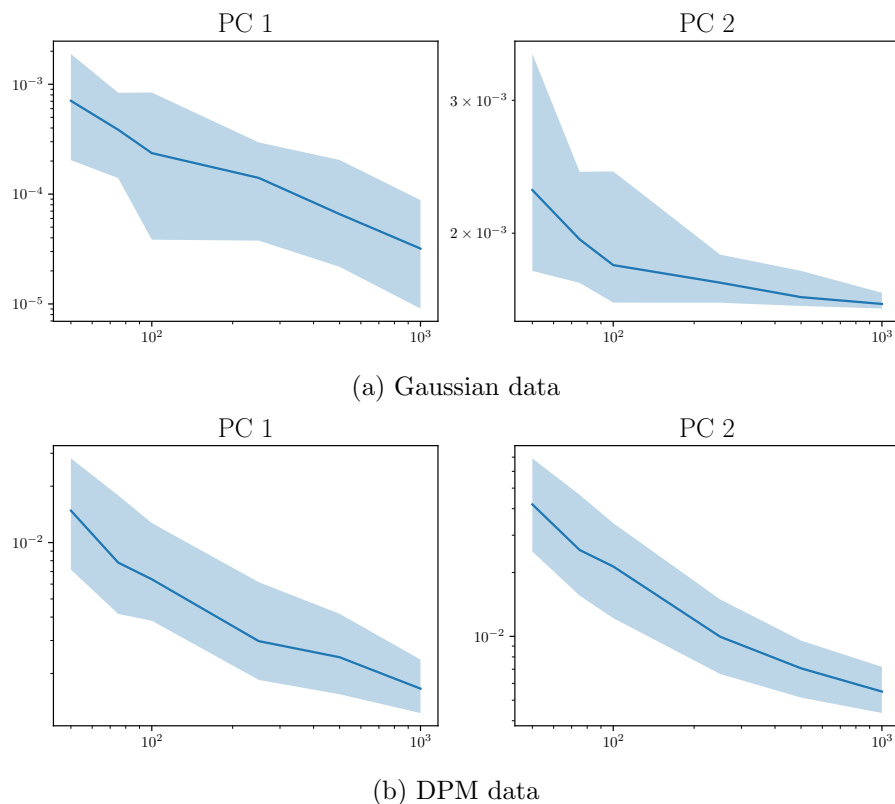
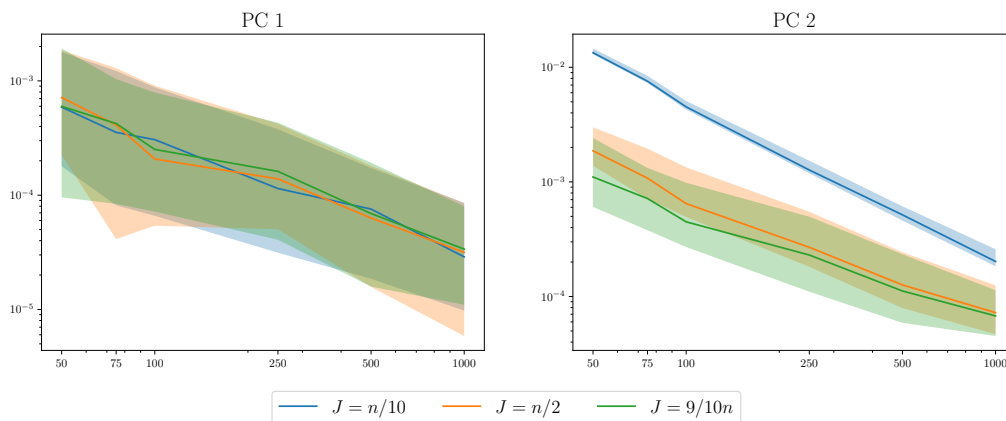


Figure 13: L_2 distance between estimated and true principal directions when $J = 20$ as a function of n . Solid line represents the median and the shaded area to a 90% confidence interval estimated from 100 independent repetition.

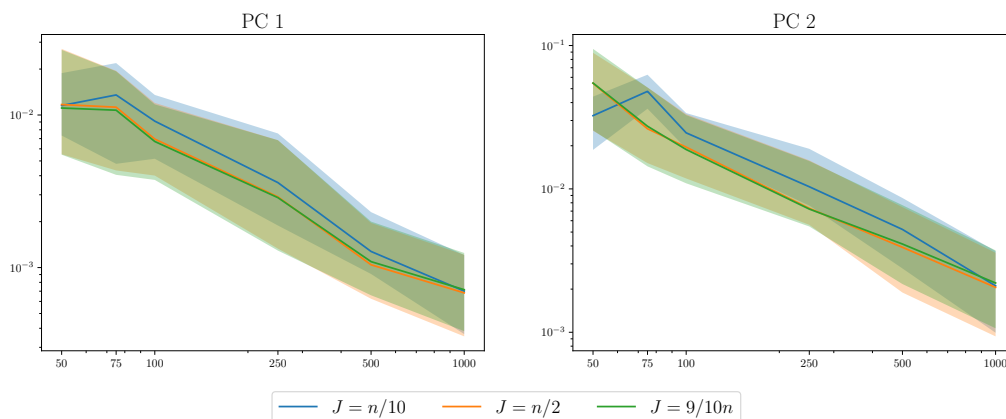
Figure 13 shows the case of fixed J for both data generation strategies. It is clear that in both cases the error quickly decreases to zero (observe that both the x and y axes are in log scale), but the convergence speed is surely sub-exponential when looking, for instance, at the second principal direction.

When increasing the number of basis elements with n , we consider three strategies letting $J = n/10$, $n/2$ and $9/10n$ respectively (rounded to the closest integer). Figure 14 shows the errors between the true and estimated principal directions in this case. Note that the convergence rate looks exponential for both data generating processes for every choice of $J = J(n)$ (increasing with n). In the case of Gaussian data, we observe smaller errors (as low as 10^{-5} for the first direction and 10^{-4} for the second direction) than in the case of the more challenging DPM data set, see Figure 14. For the former data set, using a large number of basis functions such as $9/10n$ or $n/2$ provides a much better fit than using $n/10$ basis functions on the second principal direction. For DPM data, the errors are in general two orders of magnitude higher than with Gaussian data. This is likely due to the different data generating process, which results in a more challenging problem. Interestingly, the errors are almost equal for all values of J (when fixing n).

Let us now analyze the projected regression. The independent variable are generated similarly to Section 8, by discretizing the interval $[0, 1]$ in 1,000 equispaced intervals, the value of the quantile function $F_{z_i}^-$ in the j -th interval equals $\sum_{k=1}^j \delta_{ik}$ and $(\delta_{i1}, \dots, \delta_{i1000}) \sim \text{Dirichlet}(0.01, \dots, 0.01) + \mathcal{U}([0, 5])$. We fix the kernel $\beta^*(t, s)$ (details are given below) and let quantile functions $F^{yi} = \Pi_{L_2([0,1]^t)} \circ \Gamma_{\beta^*}(F_{z_i}^-) + \mathcal{N}(0, (0.1)^2)$.



(a) Gaussian data



(b) DPM data

Figure 14: L_2 distance between estimated and true principal directions as a function of n for different choices of J . Solid line represents the median and the shaded area to a 90% confidence interval estimated from 100 independent repetition.

We consider two different choices of β^* : a smooth function $\beta_1^*(t, s) = (t - 1/2)^3 + (s - 1/2)^3$, for which we expect that a small number of spline basis will give a low error, and a rougher function $\beta_2^*(t, s)$ defined as

$$\beta_2^*(t, s) = \sum_{k,h=1}^{10} \beta_1^*(0.1k, 0.1h) \mathbb{I}[(t, s) \in [0.1(k-1), 0.1k] \times [0.1(h-1), 0.1h]]$$

that is, β_2^* corresponds to an approximation of β_1^* on a 10×10 grid. As in the case of PCA, we present two simulations for each choice of β_i^* , $i=1,2$, where we first fix the number of spline basis $J = 20$ while increasing the sample size n and second compare the performance for various values of J . We do not adopt the same strategy of setting J as a fraction of the number of n since the number of parameters to estimates grows quadratically with J which makes the computational cost substantial when $J \geq 100$. We measure both the seminorm error $\|\hat{\beta} - \beta^*\|_{C_Z}$ and the mean square prediction error on an unseen “test” set of 1,000 samples.

Figure 15 shows the seminorm error and the prediction error when $J = 20$ as n increases, while in Figure 16 various values of J are also considered. When data are generated from β_1^* , $J = 20$ spline basis is more than enough (and actually $J = 10$ would suffice) and

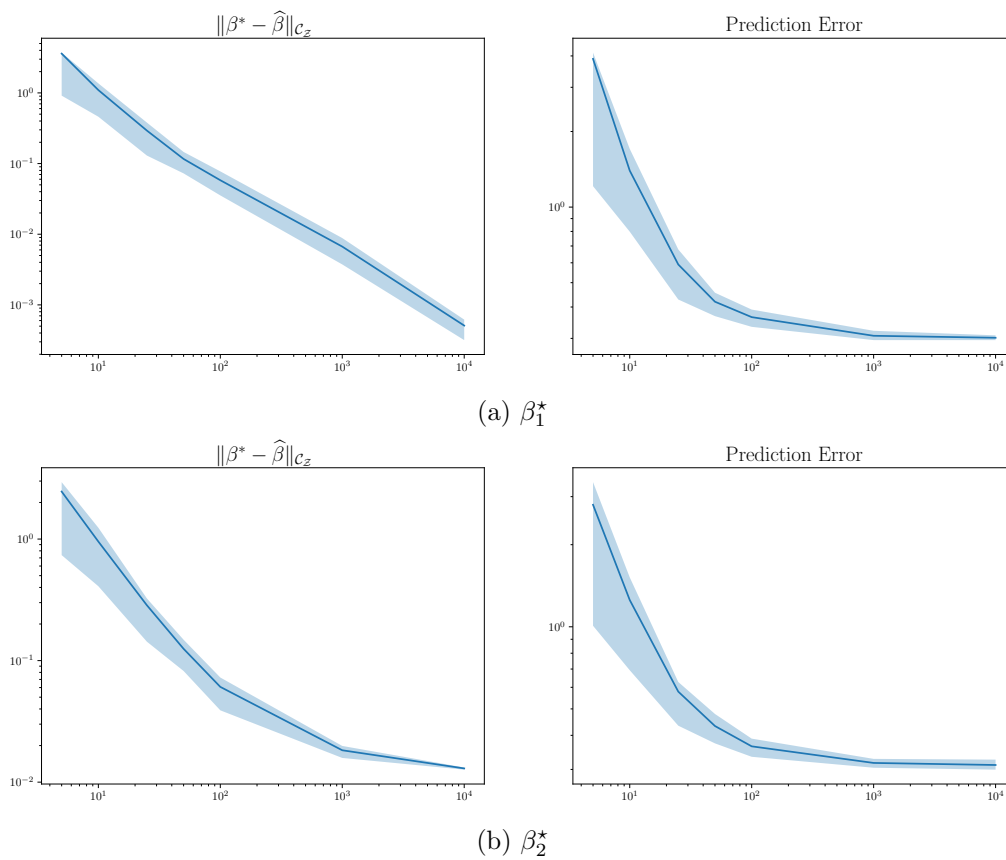


Figure 15: Seminorm error (left) and mean square prediction error (right) for different choices of the kernel used to generate data, when $J = 20$ as a function of n . Solid line represents the median and the shaded area to a 90% confidence interval estimated from 100 independent repetition.

the seminorm error in Figure 15(a) and Figure 16(a) decays exponentially while the prediction error reaches the irreducible error with $n = 10^3$ samples. When data are generated from β_2^* the seminorm error does not show the same exponential decay when $J = 20$ (see Figure 15(b)), but it does for larger values of J , in particular it seems that the error obtained with $J = 50$ is the same obtained when $J = 100$, see Figure 16(b). Hence, it is clear that the choice of J is crucial to obtain a fast decay of the error: when the kernel to be approximated is not very smooth, a larger values of spline basis elements are needed, as one would expect.

We conclude this discussion by giving a practical advice on how to select J for a given data set. Our suggestion is to let J to be the smallest value that allows for a reconstruction error smaller than a given threshold, which may depend on the specific inferential task. For instance, if the problem is PCA and the goal is to provide a descriptive analysis of the variability, a (relative) approximation error below 0.05 will typically give satisfactory results. If instead the goal is only to perform dimensionality reduction and working on the scores of a PCA as Euclidean data, one should aim for a lower approximation error, possibly of the order of 10^{-4} . A similar reasoning can be applied to the regression: if the goal is mainly to interpret the estimate $\hat{\beta}$ a larger reconstruction error can be allowed. If instead one is interested in obtaining very accurate predictions, a lower error is preferred. For instance, when β_1^* is used to generate the data, the reconstruction error for both

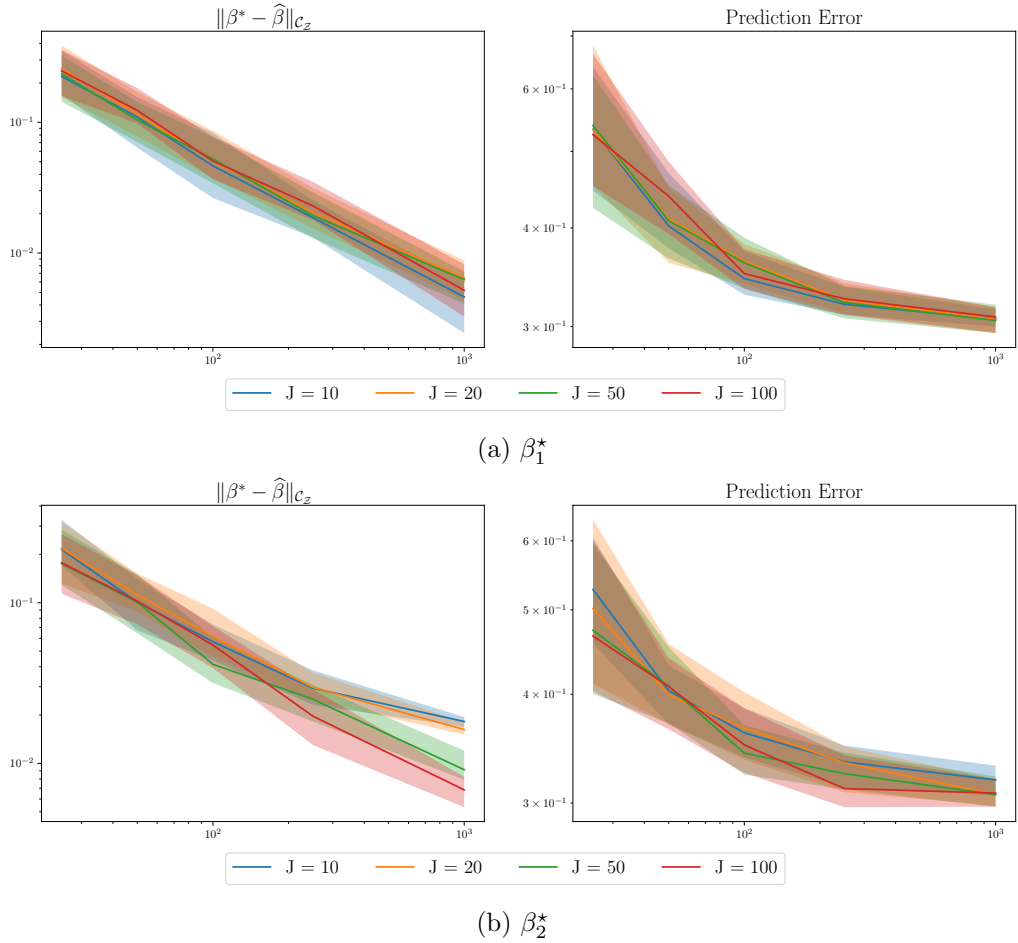


Figure 16: Seminorm error (left) and mean square prediction error (right) for different choices of the kernel used to generate data, as a function of n for different values of J . Solid line represents the median and the shaded area to a 90% confidence interval estimated from 100 independent repetition.

dependent and independent variables is below 10^{-4} for $J \geq 20$, while to get to the same error when β_2^* is used one must use $J = 100$ basis.

Appendix D. Limitations of the projected framework

D.1 When the projected PCA performs poorly

Here, we show an example to highlight the limitations of the proposed framework, specifically of the projected PCA. The main idea behind this example is that the projected principal directions will be different from the nested geodesic ones when data are concentrated around the “borders” of X , as in the trivial example shown in Figure 1. In the Wasserstein case, X is the space of quantile functions so that the border composed of functions that are constant on a subset of $[0, 1]$.

Hence, we consider the following data generating process, modeling directly the quantile functions

$$F_i^-(t) = \begin{cases} v_{i1}, & \text{if } t < 0.5 \\ v_{i1} + v_{i2}, & \text{if } t > 0.5 \end{cases}$$

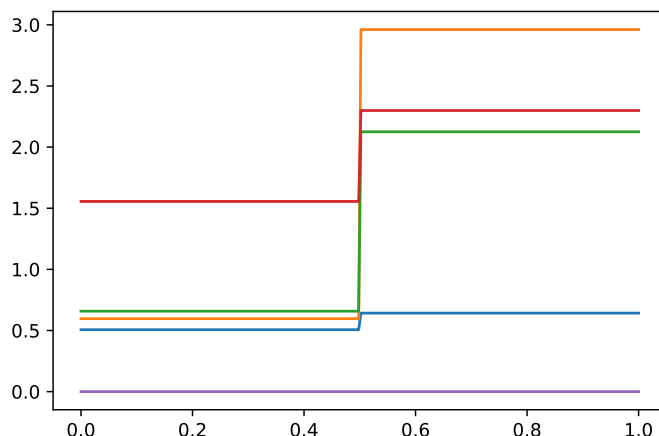


Figure 17: Five quantile functions from the data generating process considered in Appendix D.1

where $v_{ij} \sim \max\{0, \mathcal{N}(0, 1)\}$ independently. See Figure 17 for a random sample from this data generating process.

In this case, computing the projected PCA results in an interpretability score IS_k equal to one for $k = 1, 2$ and equal to zero for $k = 3, 4, \dots$. Hence, from the third principal direction onward, the projected PCA does not give any reliable information and, if those directions are needed, in this case a nested PCA could be preferred. Despite the poor interpretability scores from the third direction onward, the reconstruction errors are always good as $NRE_1 = 0.26$ and $NRE_k \approx 10^{-6}$ for $k \geq 2$. Moreover, the ghost variances GV_k are smaller than 10^{-10} for all values of k , so that this particular data set would be a good candidate for the hybrid methods mentioned in Section 7.2.

In summary, in our experience, the performance of the projected PCA can suffer when considering the interpretability of the directions associated to lower variability, but usually (at least always in our examples) gives a reasonable reconstruction error and ghost variance.

D.2 Inconsistent scores when increasing dimensions

Here, we highlight a feature which is shared by both projected and nested PCA, that is, the scores of the projection onto a projected principal component are dependent on the dimension of the principal component, as already noted in Section 3.1.

This can be considered a limitation to those frameworks, because it contributes to the complexity of the analysis: one has always to fix the dimension of the chosen principal component and use the scores accordingly obtained. For instance, the scores, both for nested and projected PCAs, coincide with the L_2 scores when the dimension of the principal components is equal to the cardinality of the spline basis J . This happens because the principal components are not linear subspaces. As a consequence also the interpretability score of a direction is dimension-dependent.

Hence, the choice of the dimension k must be carried out balancing (i) a parsimonious representation, (ii) a low reconstruction error, so that the projections on the principal components yield good approximations of the data, and (iii) the interpretability score of the directions.

Thus, opposed to standard Euclidean PCA, where the $k+1$ -th direction does not change the behavior of the data along the previous k directions (i.e., the scores), when doing (any)

PCA in Wasserstein space the whole picture must always be taken into account, both for nested and projected PCA to assess the interpretability of the results.

Finally, note that such interpretability might be low for both intrinsic and extrinsic methods, but this means that the Wasserstein metric may not be the most adequate to capture and explain the variability of the data set.

References

- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- Dragi Anevski and Philippe Soulier. Monotone spectral density estimation. *The Annals of Statistics*, 39(1):418–438, 2011.
- Dragi Anevski, Ola Hössjer, et al. A general asymptotic scheme for inference under order restrictions. *The Annals of Statistics*, 34(4):1874–1930, 2006.
- Jean-Pierre Aubin and Hélène Frankowska. *Set-valued analysis*. Springer Science & Business Media, 2009.
- Miriam Ayer, H. D. Brunk, G. M. Ewing, W. T. Reid, and Edward Silverman. An empirical distribution function for sampling with incomplete information. *The Annals of Mathematical Statistics*, 26(4):641–647, 1955. ISSN 00034851. URL <http://www.jstor.org/stable/2236377>.
- Monami Banerjee, Rudransh Chakraborty, Edward Ofori, David Vaillancourt, and Baba C Vemuri. Nonlinear regression on riemannian manifolds and its applications to neuro-image analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 719–727. Springer, 2015.
- Federico Bassetti, Antonella Bodini, and Eugenio Regazzini. On minimum kantorovich distance estimators. *Statistics & probability letters*, 76(12):1298–1302, 2006.
- Sergey Berezin and Azat Miftakhov. On barycenters of probability measures. *arXiv preprint arXiv:1911.07680*, 2019.
- Espen Bernton, Pierre E Jacob, Mathieu Gerber, and Christian P Robert. On parameter estimation with the wasserstein distance. *Information and Inference: A Journal of the IMA*, 8(4):657–676, 2019a.
- Espen Bernton, Pierre E Jacob, Mathieu Gerber, Christian P Robert, et al. Approximate bayesian computation with the wasserstein distance. *Journal of the Royal Statistical Society Series B*, 81(2):235–269, 2019b.
- Michael J Best and Nilotpal Chakravarti. Active set algorithms for isotonic regression; a unifying framework. *Mathematical Programming*, 47(1-3):425–439, 1990.
- Rabindra N Bhattacharya, L Ellingson, X Liu, V Patrangenaru, and M Crane. Extrinsic analysis on manifolds is computationally faster than intrinsic analysis with applications to quality control by machine vision. *Applied Stochastic Models in Business and Industry*, 28(3):222–235, 2012.
- Jérémie Bigot, Raúl Gouet, Thierry Klein, Alfredo López, et al. Geodesic PCA in the Wasserstein space by convex PCA. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 53, pages 1–26. Institut Henri Poincaré, 2017.

- T. Tony Cai and Peter Hall. Prediction in functional linear regression. *The Annals of Statistics*, 34:2159–2179, 2006.
- Jiezhong Cao, Langyuan Mo, Yifan Zhang, Kui Jia, Chunhua Shen, and Mingkui Tan. Multi-marginal Wasserstein GAN. In *Advances in Neural Information Processing Systems*, pages 1776–1786, 2019.
- Marta Catalano, Antonio Lijoi, and Igor Prünster. Measuring dependence in the wasserstein distance for bayesian nonparametric models. *The Annals of Statistics*, forthcoming, 2021.
- Elsa Cazelles, Vivien Seguy, Jérémie Bigot, Marco Cuturi, and Nicolas Papadakis. Geodesic PCA versus log-PCA of histograms in the Wasserstein space. *SIAM Journal on Scientific Computing*, 40(2):B429–B456, 2018.
- Yaqing Chen, Zhenhua Lin, and Hans-Georg Müller. Wasserstein regression*. *Journal of the American Statistical Association*, 0(ja):1–40, 2021. doi: 10.1080/01621459.2021.1956937. URL <https://doi.org/10.1080/01621459.2021.1956937>.
- Marco Cuturi. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. In *Advances in Neural Information Processing Systems*, pages 2292–2300, 2013.
- Marco Cuturi and Arnaud Doucet. Fast Computation of Wasserstein Barycenters. In *International Conference on Machine Learning*, pages 685–693, 2014.
- Marco Cuturi, Olivier Teboul, and Jean-Philippe Vert. Differentiable Ranking and Sorting using Optimal Transport. In *Advances in Neural Information Processing Systems*, pages 6861–6871, 2019.
- Priyam Das and Subhashis Ghosal. Bayesian quantile regression using random B-spline series prior. *Computational Statistics & Data Analysis*, 109:121–143, 2017.
- J. Dauxois, A. Pousse, and Y. Romain. Asymptotic Theory for the Principal Component Analysis of a Vector Random Function: Some Applications to Statistical Inference. *Journal of Multivariate Analysis*, 12:136–154, 1982.
- Carl De Boor and James W Daniel. Splines with Nonnegative B-spline Coefficients. *Mathematics of computation*, 28(126):565–568, 1974.
- Pedro Delicado. Dimensionality reduction when data are density functions. *Computational Statistics & Data Analysis*, 55:401–420, 01 2011.
- Frank Deutsch. *Best Approximation in Inner-Product Spaces*. Springer Science & Business Media, 2012.
- Richard Dykstra, Tim Robertson, and Farrol T Wright. *Advances in Order Restricted Statistical Inference: Proceedings of the Symposium on Order Restricted Statistical Inference Held in Iowa City, Iowa, September 11–13, 1985*, volume 37. Springer Science & Business Media, 2012.
- Juan José Egozcue, José Luis Díaz-Barrero, and Vera Pawlowsky-Glahn. Hilbert Space of Probability Density Functions Based on Aitchison Geometry. *Acta Mathematica Sinica*, 22(4):1175–1182, 2006.

- Thomas S Ferguson. Bayesian density estimation by mixtures of normal distributions. In *Recent advances in statistics*, pages 287–302. Elsevier, 1983.
- P. Fletcher. Geodesic Regression and the Theory of Least Squares on Riemannian Manifolds. *International Journal of Computer Vision*, 105, 11 2013.
- P Thomas Fletcher, Conglin Lu, Stephen M Pizer, and Sarang Joshi. Principal geodesic analysis for the study of nonlinear statistics of shape. *IEEE transactions on medical imaging*, 23(8):995–1005, 2004.
- K Hron, Alessandra Menafoglio, Matthias Templ, Klára Hrušová, and P Filzmoser. Simplified principal component analysis for density functions in Bayes spaces. *Computational Statistics & Data Analysis*, 94:330–350, 07 2014.
- Stephan Huckemann, Thomas Hotzand, and Axel Munk. Intrinsic shape analysis: Geodesic PCA for Riemannian manifolds modulo isometric lie group actions. *Statistica Sinica*, 20:1–58, 2010.
- Stephan F Huckemann and Benjamin Eltzner. Backward nested descriptors asymptotics with inference on stem cell differentiation. *The Annals of Statistics*, 46(5):1994–2019, 2018.
- Sungkyu Jung, Ian L Dryden, and James Stephen Marron. Analysis of principal nested spheres. *Biometrika*, 99(3):551–568, 2012.
- Alois Kneip and Klaus J. Utikal. Inference for Density Families Using Functional Principal Component Analysis. *Journal of the American Statistical Association*, 96(454):519–542, 2001.
- J. Le-Rademacher and L. Billard. Principal component analysis for histogram-valued data. *Advances in Data Analysis and Classification*, 11(2):327–351, 2017.
- Apoorva Mandavilli. Why does the coronavirus hit men harder? a new clue. 08 2020. URL <https://www.nytimes.com/2020/08/26/health/coronavirus-men-immune.html>.
- Alessandra Menafoglio, Alberto Guadagnini, and Piercesare Secchi. A kriging approach based on Aitchison geometry for the characterization of particle-size curves in heterogeneous aquifers. *Stochastic Environmental Research and Risk Assessment*, 28(7):1835–1851, 2014.
- Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Histoire de l’Académie Royale des Sciences de Paris*, 1781.
- P. Nagabhushan and R. Pradeep Kumar. Histogram PCA. In Derong Liu, Shumin Fei, Zengguang Hou, Huaguang Zhang, and Changyin Sun, editors, *Advances in Neural Networks – ISNN 2007*, pages 1012–1021, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg.
- Victor M Panaretos and Yoav Zemel. *An Invitation to Statistics in Wasserstein Space*. Springer Nature, 2020.
- Vic Patrangenaru and Leif Ellingson. *Nonparametric Statistics on Manifolds and Their Application to Object Data Analysis*. CRC Press, 2015.

- Xavier Pennec. Intrinsic Statistics on Riemannian Manifolds: Basic Tools for Geometric Measurements. *Journal of Mathematical Imaging and Vision*, 25:127–154, 07 2006.
- Xavier Pennec. Statistical Computing on Manifolds: From Riemannian geometry to Computational Anatomy. In *LIX Fall Colloquium on Emerging Trends in Visual Computing*, pages 347–386. Springer, 2008.
- Xavier Pennec. Barycentric subspace analysis on manifolds. *The Annals of Statistics*, 46(6A):2711–2746, 2018.
- Gabriel Peyré, Marco Cuturi, et al. Computational Optimal Transport: With Applications to Data Science. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- Florian A Potra and Stephen J Wright. Interior-point methods. *Journal of Computational and Applied Mathematics*, 124(1-2):281–302, 2000.
- Simon Potter, Marco Del Negro, Giorgio Topa, and Wilbert Van der Klaauw. The advantages of probabilistic survey questions. *Review of Economic Analysis*, 9(1):1–32, 2017.
- Luboš Prchal and Pascal Sarda. Spline estimator for functional linear regression with functional response. *Technical Report*, 2007.
- Natalya Pya and Simon N Wood. Shape constrained additive models. *Statistics and Computing*, 25(3):543–559, 2015.
- Xin Qi and Hongyu Zhao. Some theoretical properties of Silverman’s method for smoothed functional principal component analysis. *Journal of Multivariate Analysis*, 102:741–767, 2011.
- James O Ramsay. Functional data analysis. *Encyclopedia of Statistical Sciences*, 4, 2004.
- R. Tyrrell Rockafellar and Roger J.-B. Wets. *Variational Analysis*. Springer Verlag, Heidelberg, Berlin, New York, 1998.
- Oldemar Rodríguez, Edwin Diday, and Suzanne Winsberg. Generalization of the Principal Components Analysis to Histogram Data. pages 12–16, 2000.
- Kazuyuki Sekitani and Yoshitsugu Yamamoto. A recursive algorithm for finding the minimum norm point in a polytope and a pair of closest points in two polytopes. *Mathematical Programming*, 61:233–249, 1993.
- Bernard W Silverman et al. Smoothed functional principal components analysis by choice of norm. *The Annals of Statistics*, 24(1):1–24, 1996.
- Sanvesh Srivastava, Volkan Cevher, Quoc Dinh, and David Dunson. Wasp: Scalable Bayes via barycenters of subset posteriors. In *Artificial Intelligence and Statistics*, pages 912–920, 2015.
- Sean J Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.
- Rosanna Verde, Antonio Irpino, and Antonio Balzanella. Dimension reduction techniques for distributional symbolic data. *IEEE transactions on cybernetics*, 46, 01 2015.

Cédric Villani. *Optimal Transport: old and new*, volume 338. Springer Science & Business Media, 2008.

A Waechter and LT Biegler. On the implementation of a primal-dual interior point filter line search algorithm for large-scale nonlinear programming. *Mathematical Programming*, 106:25–56, 2006.

Chao Zhang, Piotr Kokoszka, and Alexander Petersen. Wasserstein autoregressive models for density time series. *arXiv preprint arXiv:2006.12640*, 2020.