



HAL
open science

Barycentric Subspace Analysis on Manifolds

Xavier Pennec

► **To cite this version:**

Xavier Pennec. Barycentric Subspace Analysis on Manifolds. *Annals of Statistics*, In press. hal-01343881v1

HAL Id: hal-01343881

<https://hal.science/hal-01343881v1>

Submitted on 11 Jul 2016 (v1), last revised 28 Sep 2017 (v2)

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

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

Barycentric Subspace Analysis on Manifolds

Xavier Pennec*

Université Côte d'Azur and Inria Sophia-Antipolis Méditerranée
Asclepios team, Inria Sophia Antipolis
2004 Route des Lucioles, BP93
F-06902 Sophia-Antipolis Cedex, France

July 10, 2016

Abstract

This paper investigates the generalization of Principal Component Analysis (PCA) to Riemannian manifolds. We first propose a new and more general type of family of subspaces in manifolds that we call barycentric subspaces. They are implicitly defined as the locus of points which are weighted means of $k + 1$ reference points. As this definition relies on points and not on tangent vectors, it can also be extended to geodesic spaces which are not Riemannian. For instance, in stratified spaces, it naturally allows principal subspaces that span several strata, which is impossible in previous generalizations of PCA. We show that barycentric subspaces locally define a submanifold of dimension k which generalizes geodesic subspaces.

Second, we rephrase PCA in Euclidean spaces as an optimization on flags of linear subspaces (a hierarchy of properly embedded linear subspaces of increasing dimension). We show that the Euclidean PCA minimizes the sum of the unexplained variance by all the subspaces of the flag, also called the Area-Under-the-Curve (AUC) criterion. Barycentric subspaces are naturally nested, allowing the construction of hierarchically nested subspaces. Optimizing the AUC criterion to optimally approximate data points with flags of affine spans in Riemannian manifolds lead to a particularly appealing generalization of PCA on manifolds, that we call Barycentric Subspaces Analysis (BSA).

MSC-class: 60D05 (Primary) 62H25, 58C06 (Secondary)

keywords: Manifold, Fréchet mean, Barycenter, Subspaces, Flag of subspaces, PCA

*xavier.pennec@inria.fr

1 Introduction

In a Euclidean space, the principal k -dimensional affine subspace of the Principal Component Analysis (PCA) procedure is equivalently defined by minimizing the variance of the residuals (the projection of the data point to the subspace) or by maximizing the explained variance within that affine subspace. This double interpretation is available through Pythagoras' theorem, which does not hold in more general manifolds. A second important observation is that principal components of different orders are nested, enabling the forward or backward construction of nested principal components.

Generalizing PCA to manifolds first requires the definition of the equivalent of affine subspaces in manifolds. For the zero-dimensional subspace, an intrinsic generalization of the mean on manifolds naturally comes into mind: the Fréchet mean is the set of global minima of the variance, as defined by Fréchet [1948] in general metric spaces. The set of local minima of the variance was named Karcher mean by Kendall [1990] after the work of Karcher [1977] on Riemannian centers of mass (see Karcher [2014] for a discussion of the naming and earlier papers). From a statistical point of view, Bhattacharya and Patrangenaru [2003, 2005] have studied in depth the asymptotic properties of the empirical Fréchet / Karcher mean.

The one-dimensional component would then quite naturally be a geodesic passing through the mean point. Higher-order components are more difficult to define. The simplest generalization is tangent PCA (tPCA), which amounts unfolding the whole distribution in the tangent space at the mean, and computing the principal components of the covariance matrix in the tangent space. The method is thus based on the maximization of the explained variance, which is consistent with the entropy maximization definition of a Gaussian on a manifold proposed by Pennec [2006]. tPCA is actually implicitly used in most statistical works on shape spaces and Riemannian manifolds because of its simplicity and efficiency. However, if tPCA is good for analyzing data which are sufficiently centered around a central value (unimodal or Gaussian-like data), it is often not sufficient for distributions which are multimodal or supported on large compact subspaces (e.g. circles or spheres).

Instead of an analysis of the covariance matrix, Fletcher et al. [2004] proposed the minimization of least squares distances to subspaces which are totally geodesic at a point. These Geodesic Subspaces (GS) are spanned by the geodesics going through a point with tangent vector restricted to belong to a linear subspace of the tangent space. They coined the procedure Principal Geodesic Analysis (PGA). However, the least-squares procedure is computationally expensive, so that they approximated it in practice with tPCA in this paper. This led to many confusions between tPCA and PGA. A real implementation of the original PGA procedure was only recently provided by Sommer et al. [2013]. PGA is allowing to build a flag (sequences of embedded subspaces) of principal geodesic subspaces consistent with a forward component analysis approach. We build components iteratively from dimension 0 (the mean point), dimension 1 (a geodesic) and higher dimensions by selecting the direction in the tangent

space at the mean that optimally reduces the square distance of data points to the geodesic subspace. In this procedure, the mean always belongs to geodesic subspaces even when they are not part of the distribution support.

To alleviate this problem, Huckemann and Ziezold [2006], and later Huckemann et al. [2010], proposed a relaxation of the requirement that the base-point of the geodesic subspace be the Fréchet mean: they start at the first order component directly with the geodesic best fitting the data, which is not necessarily going through the mean. The second principal geodesic is chosen orthogonally to the first one, and higher order components are added orthogonally at the crossing point of the first two components. The method was named Geodesic PCA (GPCA). Further relaxing the assumption that second and higher order components should cross at a single point, Sommer [2013] proposed a parallel transport of the second direction along the first principal geodesic to define the second coordinates, and iteratively define higher order coordinates through horizontal development along the previous modes. Other principal decompositions have also been proposed, like the principal graphs of Gorban and Zinovyev [2009], extending the idea of principal points and k-means.

All the above-cited methods except the last one are intrinsically forward methods that build successively larger and larger approximation spaces for the data. A notable exception is the concept of Principal Nested Spheres (PNS), proposed by Jung et al. [2012]: a framework for non-geodesic decomposition of high-dimensional spheres used in the context of planar landmarks shape spaces. Here, subspheres are viewed as slices of a higher dimensional sphere by affine hyperplanes. In this process, the nested subsphere is not of radius one, unless the hyperplane passes through the origin. The backward analysis approach determines a decreasing family of subspace. Damon and Marron [2013] have recently generalized this approach to manifolds with the help of a nested sequence of relations. However, up to now, such a sequence of relationships was only known for spheres or Euclidean spaces.

We first propose in this paper new types of family of subspaces in manifolds: barycentric subspaces (BS). Barycentric subspaces generalize geodesic subspaces and nested spheres and can naturally be nested, allowing the construction of inductive forward or backward nested subspaces. We then rephrase PCA in Euclidean spaces as an optimization on flags of linear subspaces (a hierarchy of properly embedded linear subspaces of increasing dimension). To that end, we propose an extension of the unexplained variance criterion that generalizes nicely to flags of barycentric subspaces in Riemannian manifolds. This leads to a particularly appealing generalization of PCA on manifolds, that we call Barycentric Subspaces Analysis (BSA).

1.1 Paper Organization

Barycentric subspaces are defined in Section 2 as the locus of points which are weighted means of $k + 1$ reference points. Depending on the generalization of the mean that we use on manifolds, Fréchet mean, Karcher mean or exponential barycenter, we obtain the Fréchet / Karcher or Exponential (FBS / KBS / EBS)

barycentric subspaces. As the definition relies on points and not explicitly on tangent vectors for geodesic parametrization, an interesting side effect is that BS can also be extended to more general geodesic spaces that are not Riemannian. For instance, in stratified spaces, barycentric subspaces may naturally span several strata. For Riemannian manifolds, we show that these definition are highly related since they are subsets of each other (except possibly at the cut locus of the reference points). The EBS is the largest of these barycentric subspaces. Its implicit definition exhibits some affine properties which do not depend on the metric. We define the affine span as the closure of the EBS. When the manifold \mathcal{M} is complete, this implies that the affine span is also complete. A first draft of these definitions was proposed in Pennec [2015] without any proof. In this preliminary work, one should be careful that the affine span had a different definition as it was just a synonym for the EBS.

In generic conditions, we show in Section 3 that the regular part of a barycentric subspace is a stratified space which is locally a submanifold of dimension k . Its singular set of dimension $k - l$ corresponds to the case where l of the reference point belongs to the affine span defined by the $k - l$ other reference points. In non-generic conditions, points may coalesce along certain directions, defining non local jets¹ instead of a regular $k + 1$ -tuple. Restricted geodesic subspaces, which are defined by k tangent vectors at a point, correspond to the limit of the affine span when the k -tuple converges towards that jet.

We exemplify in Section 4 the equations of barycentric subspaces in one of the simplest manifold: the sphere. We show that the affine span of $k + 1$ different reference points on the n -dimensional sphere is the k -dimensional great subsphere that contains the reference points. In fact, any $k + 1$ -tuple of points of that great k -dimensional subsphere generates the same affine span, which is also a geodesic subspace. This coincidence of spaces is due to the very high symmetry of the sphere. For second order jets, we show that we obtain subspheres of different radii, which show that Principal Nested Spheres are also a limit case of affine spans. We conjecture that this can be generalized to higher order derivatives in general manifolds using techniques from sub-Riemannian geometry. This way, some non-geodesic decomposition schemes such as loxodromes and splines could probably also be seen as limit cases of barycentric subspaces. Determining which of the points of the spherical affine span belong to the Karcher barycentric subspaces (KBS) turns out to be a difficult algebraic problem. Simple numerical tests with random data show that the index of the Hessian of the variance at critical points can be arbitrary, thus subdividing the EBS into many regions. As a result, the KBS covers only a small portion of the subsphere containing the reference points in generic conditions. This suggests that the affine span might be a much more interesting definition for subspace analysis purposes.

Finally, we discuss in Section 5 the use of these barycentric subspaces to generalize PCA on manifolds. BS can be naturally nested by defining an or-

¹Non-local jets, or multijets, generalize subspaces of the tangent spaces to higher differential orders with multiple base points.

dering of the reference points. Like for PGA, this enables the construction of a forward nested sequence of subspaces which contains the Fréchet mean. In addition, BSA also provides backward nested sequences which may not contain the mean. However, the criterion on which these constructions are based can be optimized for each subspace independently but not consistently for the whole sequence of subspaces. In order to obtain a global criterion, we rephrase PCA in Euclidean spaces as an optimization on flags of linear subspaces (a hierarchies of properly embedded linear subspaces of increasing dimension). To that end, we propose an extension of the unexplained variance criterion (the Area-Under-the-Curve criterion) that generalizes nicely to flags of affine spans in Riemannian manifolds. This results into a particularly appealing generalization of PCA on manifolds, that we call Barycentric Subspaces Analysis (BSA).

1.2 Riemannian geometry

We summarize in this section the notations used for Riemannian manifolds. A more detailed introduction to these notations can be found in Penneec [2006] and in A. We consider a differential manifold \mathcal{M} provided with a smooth scalar products $\langle \cdot | \cdot \rangle_x$ called the Riemannian metric on each tangent space $T_x\mathcal{M}$ at point x of \mathcal{M} . In a chart, the metric is fully specified by the dot product of the tangent vector to the coordinate curves: $g_{ij}(x) = \langle \partial_i | \partial_j \rangle$. The Riemannian distance between any two points on \mathcal{M} is the infimum of the length of the curves joining these points. Geodesics are defined as the critical points of the energy functional. Geodesics are parametrized by arc-length in addition to optimizing the length functional. We assume in this paper that the manifold is geodesically complete, i.e. that the definition domain of all geodesics can be extended to \mathbb{R} . This means that the manifold has no boundary nor any singular point that we can reach in a finite time. As an important consequence, the Hopf-Rinow-De Rham theorem states that there always exists at least one minimizing geodesic between any two points of the manifold (i.e. whose length is the distance between the two points).

Normal coordinate system From the theory of second order differential equations, we know that there exists one and only one geodesic $\gamma_{(x,v)}(t)$ starting from the point x with the tangent vector $v \in T_x\mathcal{M}$. The exponential map at point x maps each tangent vector $v \in T_x\mathcal{M}$ to the point of the manifold that is reached after a unit time by the geodesic: $\exp_x(v) = \gamma_{(x,v)}(1)$. The exponential map is locally one-to-one around 0: we denote by $\overrightarrow{xy} = \log_x(y)$ its inverse. The maximal domain $D(x) \subset T_x\mathcal{M}$ containing 0 where the exponential map is a diffeomorphism is a connected star-shape domain limited by the tangential cut locus $\partial D(x) = C(x) \subset T_x\mathcal{M}$ (the set of vectors tv where the geodesic $\gamma_{(x,v)}(t)$ ceases to be length minimizing). Its image by the exponential map is the cut locus $\mathcal{C}(x) = \exp_x(C(x)) \subset \mathcal{M}$. This is the closure of the set of points where several minimizing geodesics starting from x meet. The image of the domain $D(x)$ by the exponential map covers all the manifold except the cut locus, which has a null measure. When the tangent space is provided with an orthonormal

basis, the chart realized by the exp and log maps is called an normal coordinate systems at x . A set of normal coordinate systems at each point of the manifold realize an atlas which is the basis of programming on Riemannian manifolds as exemplified in Pennec et al. [2006].

Differential of the Riemannian log On $\mathcal{M}/C(y)$, the gradient of the squared distance $d_y^2(x) = \text{dist}^2(x, y)$ with respect to the fixed point y is well defined and is equal to $\nabla d_y^2(x) = -2\log_x(x_i)$. The Hessian operator of the square distance is thus directly related to the differential of the log map: $\nabla^2 d_y^2(x) = -2(D_x \log_x(x_i))$. It can also be written in terms of derivatives of the exponential map as $\nabla^2 d_y^2(x) = (D \exp_x|_{\vec{x}\vec{y}})^{-1} D_x \exp_x|_{\vec{x}\vec{y}}$ to make more explicit the link with Jacobi fields. Following Brewin [2009], we computed in A the Taylor expansion of this matrix in a normal coordinate system at x :

$$- [D_x \log_x(y)]_b^a = \delta_b^a - \frac{1}{3} R_{cbd}^a \vec{x}\vec{y}^c \vec{x}\vec{y}^d - \frac{1}{12} \nabla_c R_{dbe}^a \vec{x}\vec{y}^c \vec{x}\vec{y}^d \vec{x}\vec{y}^e + O(\epsilon^3). \quad (1)$$

Here, $R_{cbd}^a(x)$ are the coefficients of the curvature tensor at x . Since we are in a normal coordinate system, the zeroth order term is the identity matrix, like in the Euclidean space, and the first order term vanishes. The Riemannian curvature tensor appear in the second order term and its covariant derivative in the third order term. It is important to see that the curvature is the leading term that makes this matrix departing from the identity (i.e. the Euclidean case) and which may lead to the non invertibility of the differential.

1.3 Moments of point distributions

In the following, we will intensively use a set of $(k + 1)$ points on a Manifold. Adding weights $(\lambda_0, \dots, \lambda_k)$ that do not sum up to zero to each point, we may see these weighted points as the sum of weighted Diracs $\mu(x) = \sum_i \lambda_i \delta_{x_i}(x)$. As this distribution is not normalized and weights can be negative, it is generally not a probability. This distribution is singular in the sense that it is not uniformly dominated by the Riemannian measure. Thus, we have to take extra care in defining its moments as the Riemannian log and distance functions are not smooth at the cut-locus of the points.

Definition 1 ($(k + 1)$ -pointed / punctured Riemannian manifold).

Let $\{x_0, \dots, x_k\} \in \mathcal{M}^{k+1}$ be a set of $k + 1$ reference points in the Riemannian manifold \mathcal{M} and $C(x_0, \dots, x_k) = \cup_{i=0}^k C(x_i)$ be the union of the cut loci of these points. We call $(k + 1)$ -pointed manifold the object consisting of the smooth manifold M and the $k + 1$ reference points, and $(k + 1)$ -punctured manifold the submanifold $\mathcal{M}^*(x_0, \dots, x_k) = \mathcal{M}/C(x_0, \dots, x_k)$ of the non-cut points of the $k + 1$ reference points.

On $\mathcal{M}^*(x_0, \dots, x_k)$, the distance to the points $\{x_0, \dots, x_k\}$ is smooth. The Riemannian log function $\vec{x}\vec{x}_i = \log_x(x_i)$ is also well defined for all the points of $\mathcal{M}^*(x_0, \dots, x_k)$ but becomes multivalued at the cut locus $C(x_i)$ of x_i . Since the

cut locus of each point is closed and has null measure, the punctured manifold $\mathcal{M}^*(x_0, \dots, x_k)$ is open and dense in \mathcal{M} , which means that it is a submanifold of \mathcal{M} . However, this submanifold is not necessarily connected. For instance in the flat torus $(S_1)^n$, the cut-locus of $k + 1$ points divides the torus into k^n disconnected cells.

Definition 2 (Weighted moments of a $(k + 1)$ -pointed manifold).
Let $(\lambda_0, \dots, \lambda_k) \in \mathbb{R}^{k+1}$ such that $\sum_i \lambda_i \neq 0$. The weighted n -order moment of a $(k + 1)$ -pointed Riemannian manifold is the n -contravariant tensor:

$$\mathfrak{M}_n(x, \lambda) = \sum_i \lambda_i \underbrace{\overrightarrow{xx_i} \otimes \overrightarrow{xx_i} \dots \otimes \overrightarrow{xx_i}}_{n \text{ times}}, \quad (2)$$

and the normalized weighted n -order moment is:

$$\underline{\mathfrak{M}}_n(x, \lambda) = \mathfrak{M}_n(x, \lambda) / \mathfrak{M}_0(\lambda). \quad (3)$$

Both tensors are smoothly defined on the punctured manifold $\mathcal{M}^*(x_0, \dots, x_k)$.

The 0-th order moment $\mathfrak{M}_0(\lambda) = \sum_i \lambda_i = \mathbf{1}^T \lambda$ is the mass. The n -th order moment is homogeneous of degree 1 in λ while the normalized n -th order moment is naturally invariant by a change of scale of the weights thanks to the use of the normalized weights. For a fixed weight λ , the first order moment $\mathfrak{M}_1(x, \lambda) = \sum_i \lambda_i \overrightarrow{xx_i}$ is a smooth vector field on the manifold $\mathcal{M}^*(x_0, \dots, x_k)$ whose zeros will be the subject of our interest. The second and higher order moments are smooth $(n, 0)$ tensor fields that will be used later through their contraction with the Riemannian curvature tensor.

2 Barycentric subspaces

2.1 Affine subspaces in a Euclidean space

In Euclidean PCA, a zero dimensional space is a point, a one-dimensional space is a line, and an affine subspace of dimension k is generated by a point and k non-collinear vectors. Alternatively, one could also generate such a subspace by taking the affine hull of $k + 1$ points in general position: $\text{Aff}(x_0, \dots, x_k) = \left\{ x = \sum_i \lambda_i x_i, \text{ with } \sum_{i=0}^k \lambda_i = 1 \right\}$. While the first definition parametrizes the space of affine subspaces with a contact element of dimension k , the second definition relies on the configuration space of $k + 1$ points in general position (a k -simplex). These two definitions are of course equivalent in a Euclidean space, but turn out to have different generalizations in manifolds. It is worth noticing that when the points are not in general conditions, the affine span is still well defined but has a lower dimensionality.

When there exists a vector of coefficients $\lambda = (\lambda_0 : \lambda_1 : \dots : \lambda_k) \in \mathbb{R}^{k+1}$ (with do not sum to zero) such that $\sum_{i=0}^k \lambda_i (x_i - x) = 0$, then λ is called the barycentric coordinates of the point x with respect to the k -simplex $\{x_0, \dots, x_k\}$.

Barycentric coordinates are not unique since they are homogeneous of degree one. Thus, one usually renormalize the coordinates by the total mass so that $\sum_{i=0}^k \lambda_i = 1$. In that case, the vertices of the simplex have the coordinates $(1, 0, 0, \dots, 0) \dots (0, 0, 0, \dots, 1)$.

Definition 3 (Projective space of barycentric coordinates (weights)). *Barycentric coordinates of $k + 1$ points live in the projective space \mathcal{P}_k minus the orthogonal of the line element $\mathbf{1} = (1 : 1 : \dots : 1)$:*

$$\mathcal{P}_k^* = \left\{ (\lambda_0 : \lambda_1 : \dots : \lambda_k) \in \mathbb{R}^{k+1} \text{ s.t. } \sum_{i=0}^k \lambda_i \neq 0 \right\}.$$

Standard charts of this space are given either by the intersection of the line elements with the "upper" unit sphere S_k of \mathbb{R}^{k+1} with north pole $\mathbf{1}/\sqrt{k}$ or by the k -plane of \mathbb{R}^{k+1} passing through the point $\mathbf{1}/k$ and orthogonal to this vector. We call normalized weights $\lambda_i = \lambda_i / (\sum_{j=0}^k \lambda_j)$ this last projection.

2.2 Fréchet / Karcher Barycentric subspaces metric spaces

The reformulation of the affine span as the weighted mean of $(k + 1)$ points for some weights suggests first to generalize the definition to metric manifolds using the Fréchet or the Karcher mean.

Definition 4 (Fréchet / Karcher barycentric subspaces of $k + 1$ points). *Let $(\mathcal{M}, \text{dist})$ be a metric space and $(x_0, \dots, x_k) \in \mathcal{M}^k$ be $k + 1$ distinct reference points. The (normalized) weighted variance at point x with weight $\lambda \in \mathcal{P}_k^*$ is:*

$$\sigma^2(x, \lambda) = \frac{1}{2} \sum_{i=0}^k \lambda_i \text{dist}^2(x, x_i) = \frac{1}{2} \sum_{i=0}^k \lambda_i \text{dist}^2(x, x_i) / \left(\sum_{j=0}^k \lambda_j \right).$$

The Fréchet barycentric subspace of these points is the locus of weighted Fréchet means of these points, i.e. the set of absolute minima of the weighted variance:

$$FBS(x_0, \dots, x_k) = \left\{ \arg \min_{x \in \mathcal{M}} \sigma^2(x, \lambda), \lambda \in \mathcal{P}_k^* \right\}$$

The Karcher barycentric subspaces $KBS(x_0, \dots, x_k)$ are defined similarly with local minima instead of global ones.

This definition restores the full symmetry between all parameters of the subspaces, contrarily to the geodesic subspaces which are intrinsically privileging one point. This definition is also sufficiently general to work on metric spaces more general than Riemannian manifolds. In stratified metric spaces, for instance, the barycentric subspace spanned by points belonging to different strata naturally maps over several strata. This is a significant improvement over geodesic subspaces used in PGA which can only be defined within a regular strata.

The reference points could be seen as landmarks in the manifold. However, since landmark has a specific meaning in morphometry, we prefer not to use this terminology. In biology, Shoval et al. [2012] used archetype points in a similar way to characterize the geometry of the phenotype space using Pareto optimality. Here, "archetype" suggests that these points are extremal in some way (for instance barycentric coordinates should be between -1 and 1), which is interesting by not mandatory in our framework.

2.3 Exponential Barycentric Subspace (EBS) and Affine Span in Riemannian manifolds

A third definition of the mean in manifolds can be used to define barycentric subspaces: exponential barycenters.

Definition 5 (Barycentric coordinates in a $(k + 1)$ -pointed manifold). *A point $x \in \mathcal{M}^*(x_0, \dots, x_k)$ has barycentric coordinates $\lambda \in \mathcal{P}_k^*$ if*

$$\mathfrak{M}_1(x, \lambda) = \sum_{i=0}^k \lambda_i \overrightarrow{xx_i} = 0. \quad (4)$$

Since the Riemannian log function $\overrightarrow{xx_i} = \log_x(x_i)$ is multiply defined on the cut locus of x_i , this definition cannot be extended to the the union of all cut loci $C(x_0, \dots, x_k)$, which is why we exclude this set and restrict the definition to $\mathcal{M}^*(x_0, \dots, x_k)$ in the present work.

Definition 6 (Exponential Barycentric Subspace (EBS)). *The EBS of the points $(x_0, \dots, x_k) \in \mathcal{M}^k$ is the set of weighted exponential barycenters of the reference points in $\mathcal{M}^*(x_0, \dots, x_k)$:*

$$EBS(x_0, \dots, x_k) = \{x \in \mathcal{M}^*(x_0, \dots, x_k) \mid \exists \lambda \in \mathcal{P}_k^* : \mathfrak{M}_1(x, \lambda) = 0\}.$$

On the punctured manifold $\mathcal{M}^*(x_0, \dots, x_k)$, the gradient of the squared distance $d_{x_i}^2(x) = \text{dist}^2(x, x_i)$ is well defined and is equal to $\nabla d_{x_i}^2(x) = -2 \log_x(x_i)$. Thus, one recognizes that Eq.(4) defines nothing else than the critical points of the variance $\sigma^2(x, \lambda) = \frac{1}{2} \sum_i \lambda_i \text{dist}^2(x, x_i)$. The EBS is thus a superset of the FBS / KBS in $\mathcal{M}^*(x_0, \dots, x_k)$.

The discontinuity of the Riemannian log on the cut locus of the reference points may hide the continuity or discontinuities of the exponential barycentric subspace. In order to ensure the completeness of the subspace and potentially reconnect different components, we define consider the closure of this set.

Definition 7 (Affine span of $(k + 1)$ points in a Riemannian manifold). *The affine span is the closure of the EBS in \mathcal{M} : $\text{Aff}(x_0, \dots, x_k) = \overline{EBS}(x_0, \dots, x_k)$. Because we assumed that \mathcal{M} is geodesically complete, this is equivalent to the metric completion of the EBS.*

The local minima of the variance which are potentially located on the cut-locus of the reference points are not part of the EBS but they are recovered

in the affine span thanks to the metric completion. FBS and KBS are thus included in the affine span, and the affine span is the largest of the barycentric subspaces.

For the following, it is interesting to introduce the dual space of admissible barycentric coordinates at each point of $x \in \mathcal{M}^*(x_0, \dots, x_k)$:

Proposition 1 (Dual subspace of admissible barycentric weights). *The space of valid barycentric weights $\Lambda(x) = \{\lambda \in \mathcal{P}_k^* | \mathfrak{M}_1(x, \lambda) = 0\}$ is either void, a point, or a linear subspace of \mathcal{P}_k^* .*

2.4 A SVD characterizations of the exponential barycentric subspace

Let $Z(x) = [\overrightarrow{xx_0}, \dots, \overrightarrow{xx_k}]$ be the smooth field of $n \times (k+1)$ matrices of vectors pointing from any point $x \in \mathcal{M}^*(x_0, \dots, x_k)$ to the reference points. We can rewrite the constraint $\sum_i \lambda_i \overrightarrow{xx_i} = 0$ in matrix form: $\mathfrak{M}_1(x, \lambda) = Z(x)\lambda = 0$, where λ is the $k+1$ vector of homogeneous coordinates λ_i .

Theorem 1 (Characterization of the exponential barycentric subspace). *Let $Z(x) = U(x).S(x).V(x)^T$ be a singular decomposition of the $n \times (k+1)$ matrix fields $Z(x) = [\overrightarrow{xx_0}, \dots, \overrightarrow{xx_k}]$ on $\mathcal{M}^*(x_0, \dots, x_k)$ (with singular values $\{s_i(x)\}_{0 \leq i \leq k}$ sorted in decreasing order). The barycentric subspace $\text{Aff}(x_0, \dots, x_k)$ is the zero level-set of the $k+1$ singular value $s_{k+1}(x)$ and the dual subspace of valid barycentric weights is spanned by the right singular vectors corresponding to the l vanishing singular values: $\Lambda(x) = \text{Span}(v_{k-l}, \dots, v_k)$ (it is void if $l = 0$).*

Proof. Since U and V are orthogonal matrices, $Z(x)\lambda = 0$ if and only if at least one singular value (necessarily the smallest one s_k) is null, and λ has to live in the corresponding right-singular space: $\Lambda(x) = \text{Ker}(Z(x))$. If we have only one zero singular value ($s_{k+1} = 0$ and $s_k > 0$), then λ is proportional to v_{k+1} . If l singular values vanish, then we have a higher dimensional linear subspace of solutions for λ . \square

The dimension of the dual space $\Lambda(x)$ is actually controlling the local dimension of the barycentric space, as we will see below.

2.5 Link between the different barycentric subspaces

In order to analyze the relationship between the Fréchet / Karcher / Exponential barycentric subspaces, we follow the seminal work of Karcher [1977]. First, the locus of local minima (i.e. Karcher mean) is a superset of the global minima (Fréchet mean). On the punctured manifolds $\mathcal{M}^*(x_0, \dots, x_k)$, the weighted variance is smooth and its critical points are the points of exponential barycentric subspace, which is also the restriction of the affine span to $\mathcal{M}^*(x_0, \dots, x_k)$. Among the critical points with a non-degenerate Hessian, local minima are characterized by a positive definite Hessian. When the Hessian is degenerate, we

cannot conclude on the local minimality without going to higher order differentials. Thus, we have

$$FBS \cap \mathcal{M}^* \subset KBS \cap \mathcal{M}^* \subset Aff \cap \mathcal{M}^* = EBS.$$

The goal of this section is to decompose the EBS into cells according to the index of the Hessian operator of the variance:

$$H(x, \lambda) = \nabla^2 \sigma^2(x, \lambda) = - \sum_{i=0}^k \lambda_i D_x \log_x(x_i). \quad (5)$$

Plugging the value of the Taylor expansion of the differential of the log of Eq.(20), we obtain the Taylor expansion:

$$[H(x, \lambda)]_b^a = \delta_b^a - \frac{1}{3} R_{cbd}^a(x) \mathfrak{M}_2^{cd}(x, \lambda) - \frac{1}{12} \nabla_c R_{dbe}^a(x) \mathfrak{M}_3^{cde}(x, \lambda) + O(\epsilon^4). \quad (6)$$

The key factor in this expression is the contraction of the Riemannian curvature with the weighted covariance tensor of the reference points. This contraction can be seen as an extension of the Ricci curvature tensor. Exactly as the Ricci curvature tensor encodes (through its metric trace, the scalar curvature) how the volume of an isotropic geodesic ball in the manifold deviates from the volume of the standard ball in a Euclidean space, the extended Ricci curvature encodes how the volume of the ellipsoid $\vec{x}\vec{y}^T \mathfrak{M}_2(x, \lambda) \vec{x} \leq \epsilon$ centered at the point x in the manifold deviates from the volume of the same ellipsoid in a Euclidean space (the tangent space).

Interestingly, in symmetric spaces (or even more generally locally symmetric affine spaces), the covariant derivative of the curvature is identically zero, which simplifies the formula above. We should be careful however that the curvature tensor still appears in higher terms. In the limit of null curvature, (e.g. for a locally Euclidean space like the torus), the Hessian matrix $H(x, \lambda)$ converges to the unit matrix, which means that it never vanishes. In general Riemannian manifolds, Equation (6) only gives a qualitative behavior but does not provide guaranties as it is a series involving higher order moments of the reference points. In order to obtain hard bounds on the spectrum of $H(x, \lambda)$, one has to investigate bounds on Jacobi fields using Riemannian comparison theorems, as is done for the proof of uniqueness of the Karcher and Fréchet means (see Karcher [1977], Kendall [1990], Le [2004], Afsari [2010], Yang [2011]).

Definition 8 (Degenerate, non-degenerate and positive points).

An exponential barycenter $x \in EBS(x_0, \dots, x_k)$ is degenerate (resp. non-degenerate or positive) if the Hessian matrix $H(x, \lambda)$ is singular (resp. definite or positive definite) for all λ in the the dual space of valid weights $\Lambda(x)$ (the right singular space of the zero singular value of $Z(x)$). The set of degenerate (resp. non-degenerate or positive) exponential barycenter is called the degenerate EBS and denoted $EBS^0(x_0, \dots, x_k)$ (resp. non-degenerate $EBS^*(x_0, \dots, x_k)$ or positive $EBS^+(x_0, \dots, x_k)$).

Positive points are obviously non-degenerate. In flat spaces (e.g. Euclidean), all the points of the punctured manifold are positive and non-degenerate. In curved manifolds, we may have degenerate points and non-degenerate but non-positive points, as we will see with the example of spheres. The definition of non-degenerate and positive points could be generalized to non-critical points (outside the affine span) by considering for instance the right singular space of the smallest singular value of $Z(x)$. However, this extended definition would depend on the metric that we choose for the space of weights and a renormalization of the weights (such as the one we will do for spheres in Section 4) can change the smallest singular value.

Theorem 2 (Karcher barycentric subspace and positive span).

$EBS^+(x_0, \dots, x_k)$ is the set of non-degenerate points of the Karcher barycentric subspace $KBS(x_0, \dots, x_k)$ on $\mathcal{M}^(x_0, \dots, x_k)$. In other words, the KBS is the positive EBS plus potentially some degenerate points of the affine span and some points of the cut locus of the reference points.*

3 Properties of the barycentric subspaces

3.1 Link with the convex hull

In a vector space, a point lies in the convex hull of a simplex if and only if its barycentric coordinates are all non-negative (meaning that they are between 0 and 1 with the unit sum constraint). In that case, the weights λ can be interpreted as a vector of probabilities. Consequently, barycentric coordinates are often thought to be related to convex hulls. However, in a general Riemannian manifold, the situation is quite different. When there are closed geodesics, the convex hull can reveal several disconnected components, unless one restrict to convex subsets of the manifolds as shown in Groisser [2004]. In metric spaces with negative curvature (CAT spaces), Weyenberg [2015] displays explicit examples of convex hulls of 3 points which are 3-dimensional rather than 2-dimensional as expected. In fact, the equivalence of barycentric subspaces with convex hulls only holds whenever the barycentric subspace is totally geodesic at each point, which happens for spheres (and probably for constant curvature spaces) but not for general Riemannian manifolds.

3.2 Barycentric simplex in a regular geodesic ball

We call barycentric simplex the subset of the FBS that has non-negative weights. It contains all the reference points, the geodesics segments between the reference points, and of course the Fréchet mean of the reference points. This is the generalization of a geodesic segment for 2 points, a triangle for 3 points, etc. The $(k - l)$ -faces of a k -simplex are the simplices defined by the barycentric subspace of $k - l + 1$ points among the $k + 1$. They are obtained by imposing the l remaining barycentric coordinates to be zero. In parallel to the writing of this paper, Weyenberg [2015] has investigated barycentric simplexes as extensions of

principal subspaces in the negatively curved metric spaces of trees in under the name Locus of Fréchet mean (LFM), with very interesting results.

Theorem 3 (Barycentric simplex in a regular geodesic ball). *Let κ be an upper bound of sectional curvatures of \mathcal{M} and $\text{inj}(\mathcal{M})$ be the radius of injection (which can be infinite) of the Riemannian manifold. Let $X = \{x_0, \dots, x_k\} \in \mathcal{M}^{(k+1)}$ be a set of points included in a regular geodesic ball $B(x, \rho)$ with $\rho < \frac{1}{2} \min\{\text{inj}(\mathcal{M}), \frac{1}{2}\pi/\sqrt{\kappa}\}$ ($\pi/\sqrt{\kappa}$ being infinite if $\kappa < 0$). The barycentric simplex is the graph of a k -dimensional differentiable function from the non-negative quadrant of homogeneous coordinates $(\mathcal{P}_k^*)^+$ to $B(x, \rho)$ and is thus at most k -dimensional. The $k-l$ -faces of the simplex are the simplices defined by the barycentric subspace of $k-l+1$ points among the $k+1$ and include the reference points themselves as vertices (0-faces) and the geodesics joining them as edges (1-faces).*

Proof. The proof closely follows the one of Karcher [1977] for the uniqueness of the Riemannian barycenter. The main argument is that $\mu_{(X, \lambda)}(x) = \sum \lambda_i \delta_{x_i}(x)$ is a probability distribution whose support is included in the strongly convex geodesic ball $B(x, \rho)$. Following Karcher [1977], the variance $\sigma^2(x, \lambda) = \frac{1}{2} \sum_i \lambda_i d^2(x, x_i)$ is strictly convex on that ball and has a unique minimum $x_\lambda \in B(x, \rho)$, which is necessarily the weighted Fréchet mean. This proof of the uniqueness of the weighted Fréchet mean with non-negative weights was actually already present in Buser and Karcher [1981]. We supplement the proof here by noting that since the Hessian $H(x_\lambda, \lambda) = \sum_i \lambda_i H_i(x_\lambda)$ is the convex combination of positive matrices, it is positive definite for all $\lambda \in (\mathcal{P}_k^*)^+$ in the positive quadrant. Thus the function x_λ is differentiable thanks to the implicit function theorem: $D_\lambda x_\lambda = H(x_\lambda, \lambda)^{(-1)} Z(x_\lambda)$. The rank of this derivative is at most k since $Z(x_\lambda) = 0$, which proves that the graph of the function x_λ describes at most a k dimensional subset in \mathcal{M} . As we will see below, it is actually a stratified space in generic conditions. \square

3.3 Local dimension of the barycentric subspaces

Let x be point of the EBS verifying $Z(x)\lambda = 0$ for some $\lambda \in \Lambda(x)$. This expression is smooth in x and λ so that we can take a Taylor expansion: at the first order, a variation of barycentric coordinates $\delta\lambda$ induces a variation of position δx which are linked through $H(x, \lambda)\delta x + Z(x)\delta\lambda = 0$. Thus, at regular points, we have

$$\delta x = -H(x, \lambda)^{(-1)} Z(x)\delta\lambda.$$

Let $Z(x) = U(x)S(x)V(x)^\top$ be a singular value decomposition with singular values sorted in decreasing order. Since x belongs to the EBS, there is at least one (say $m \geq 1$) singular value that vanish and the dual space of admissible weights is $\Lambda(x) = \text{Span}(v_{k-m}, \dots, v_k)$. For a variation of weights $\delta\lambda$ in this subspace, there is no change of coordinates, while any variation of weights in $\text{Span}(v_0, \dots, v_{k-m-1})$ induces a non-zero position variation. Thus, the tangent space of the EBS restricts to the $(k-m)$ -dimensional linear space generated by

$\{\delta x'_i = -H(x, \lambda)^{(-1)} u_i\}_{0 \leq i \leq k-m}$. Here, we see that the Hessian matrix $H(x, \lambda)$ encodes the distortion of the orthonormal frame fields $u_1(x), \dots, u_k(x)$ to match the tangent space. Since the lower dimensional subspaces are included one the larger ones, we have a stratification of our k -dimensional submanifold into $k-1, k-2, \dots, 0$ -dimensional boundaries.

Theorem 4 (Dimension of the exponential barycentric subspace at non-degenerate points). *$EBS^*(x_0, \dots, x_k)$ is a stratified space of dimension k on $\mathcal{M}^*(x_0, \dots, x_k)$. On the m -dimensional strata, $Z(x)$ has exactly $k-m+1$ vanishing singular values.*

At degenerate points, $H(x, \lambda)$ is not invertible and vectors living in its kernel are also authorized, which potentially raises the dimensionality of the tangent space, even if they do not change the barycentric coordinates. Thus pathologies may appear at degenerate points in barycentric subspaces of general manifolds. In practice, this is not the case for the sphere, as we will see in the sequel, and we conjecture that this is also not the case for symmetric spaces.

3.4 Two alternative characterizations in the spirit of PCA

Theorem 5. *Let $\Omega(x) = Z(x)^T G(x) Z(x)$ be the smooth $(k+1) \times (k+1)$ matrix field on $\mathcal{M}^*(x_0, \dots, x_k)$ with components $\Omega_{ij}(x) = \langle \overrightarrow{xx_i} \mid \overrightarrow{xx_j} \rangle_x$ and $\Sigma(x) = \mathfrak{M}_2(x, \mathbf{1}) = \sum_{i=0}^k \overrightarrow{xx_i} \overrightarrow{xx_i}^T = Z(x) Z(x)^T$ be the (scaled) $n \times n$ covariance matrix field of the reference points. $EBS(x_0, \dots, x_k)$ is the zero level-set of: $\det(\Omega(x))$, the minimal eigenvalue σ_{k+1}^2 of $\Omega(x)$, the $k+1$ eigenvalue (in decreasing order) of the covariance $\Sigma(x)$.*

Proof. The constraint $\mathfrak{M}_1(x, \lambda) = 0$ is satisfied if and only if its squared norm is zero:

$$\|\mathfrak{M}_1(x, \lambda)\|_x^2 = \|\sum_i \lambda_i \overrightarrow{xx_i}\|_x^2 = \lambda^T \cdot \Omega(x) \cdot \lambda.$$

As the function is homogeneous in λ , we can restrict to unit vectors. Adding this constrains with a Lagrange multiplier to the cost function, we end-up with the Lagrangian

$$\mathcal{L}(x, \lambda, \alpha) = \lambda^T \cdot \Omega(x) \cdot \lambda + \alpha(\lambda^T \lambda - 1) \quad (7)$$

The minimum with respect to λ is obtained for the eigenvector $\mu_{k+1}(x)$ associated to the smallest eigenvalue $\sigma_{k+1}(x)$ of $\Omega(x)$ (assuming that eigenvalues are sorted in decreasing order) and we have $\|\mathfrak{M}_1(x, \mu_{k+1}(x))\|_2^2 = \sigma_{k+1}(x)$, which is null if and only if the minimal eigenvalue is zero. Thus, the barycentric subspace of $k+1$ points is the locus of rank deficient matrices $\Omega(x)$:

$$EBS(x_0, \dots, x_k) = \phi^{(-1)}(0) \quad \text{where} \quad \phi(x) = \det(\Omega(x)).$$

One may want to relate the singular values of $Z(x)$ to the eigenvalues of $\Omega(x)$. The later are the square of the singular values of $G(x)^{1/2} Z(x)$. However, the left multiplication by the square root of the metric (a non singular but non orthogonal matrix) obviously changes the singular values in general. There is

however a special case where some singular values are equal: this is for vanishing ones. The (right) kernels of $G(x)^{1/2}Z(x)$ and $Z(x)$ are indeed the same. This shows that the EBS is an affine notion rather than a metric one, contrarily to the Fréchet / Karcher barycentric subspace.

To draw the link with the $n \times n$ covariance matrix of the reference points (we intentionally dropped the usual normalization factor $1/k + 1$ to simplify the notations), let us notice first that the definition does not assume that the coordinate system is orthonormal. Thus, the eigenvalues of the covariance matrix are depending on the chosen coordinate system, unless they vanish. In fact only the joint eigenvalues of $\Sigma(x)$ and $G(x)$ really make sense, which is why this last decomposition is sometimes called the proper orthogonal decomposition (POD). Now, the singular values of $Z(x) = U(x)S(x)V(x)^T$ are also the square root of the first $k + 1$ eigenvalues of $\Sigma(x) = U(x)S^2(x)U(x)^T$, the remaining $n - k - 1$ eigenvalues being null. Similarly, the singular values of $G(x)^{1/2}Z(x)$ are the square root of the first $k + 1$ joint eigenvalues of $\Sigma(x)$ and $G(x)$. Thus, our barycentric subspace may also be characterized as the zero level-set of the $k + 1$ eigenvalue (sorted in decreasing order) of Σ (or of the joint eigenvalue of $\Sigma(x)$ and $G(x)$), and this characterization is once again independent of the basis chosen. \square

3.5 Stability of the affine span with respect to the metric power

The Fréchet (resp. Karcher) mean can be further generalized by taking a power α of the metric to define the α -variance $\sigma^\alpha(x) = \frac{1}{\alpha} \sum_{i=0}^k \text{dist}^\alpha(x, x_i)$. The absolute (resp. local) minima of this α -variance defines the median for $\alpha = 1$ and the modes for $\alpha \rightarrow 0$. This suggest that we could further generalize barycentric subspaces by taking the locus of the minima of the weighted α -variance $\sigma^\alpha(x, \lambda) = \frac{1}{\alpha} \sum_{i=0}^k \lambda_i \text{dist}^\alpha(x, x_i)$. In fact, it turns out that all these " α -subspaces" are necessarily included in the affine span, which shows this notion is really central. To see that, let us we compute the gradient of the α -variance at any point of $\mathcal{M}^*(x_0, \dots, x_k)$:

$$\nabla_x \frac{1}{\alpha} \sigma^\alpha(x, \lambda) = - \sum_{i=0}^k \lambda_i \text{dist}^{\alpha-2}(x, x_i) \log_x(x_i).$$

We see that the critical points satisfy the equation $\sum_{i=0}^k \lambda'_i \log_x(x_i) = 0$ for the new weights $\lambda'_i = \lambda_i \text{dist}^{\alpha-2}(x, x_i)$. Thus the critical points of the α -variance are simply elements of the EBS and changing the power of the metric just amounts to a reparametrization of the barycentric weights.

3.6 Restricted geodesic submanifolds are limit of affine spans

So far, we have considered that the reference points $\{x_0, \dots, x_k\}$ were distinct in order to have a chance to generate a k -dimensional subspace. Let us investigate

what is happening when all the points $\{x_i = \exp_{x_0}(\eta w_i)\}_{1 \leq i \leq k}$ are converging to x_0 at first order along k independent vectors $\{w_i\}_{1 \leq i \leq k}$. Here, we fix $w_0 = 0$ to simplify the derivations, but the proof can be easily extended to $w_0 \neq 0$ with a suitable change of coordinate system provided that $\sum_{i=0}^k w_i = 0$.

By analogy with Euclidean spaces, where a point of the affine span $y = \sum_{i=0}^k \lambda_i x_i$ may be rewritten as the point $y = x + \eta \sum_{i=1}^k \lambda_i w_i$ of the "geodesic subspace" generated by the family of vectors $\{w_i\}_{1 \leq i \leq k}$, we expect the exponential barycentric subspace $EBS(x_0, \exp_{x_0}(\eta w_1) \dots \exp_{x_0}(\eta w_k))$ to be close to the geodesic subspace

$$GS(x, w_1, \dots, w_k) = \left\{ \exp_x \left(\sum_{i=1}^k \alpha_i w_i \right) \in \mathcal{M} \text{ for } \alpha \in \mathbb{R}^k \right\}$$

generated by the k independent vectors w_1, \dots, w_k at x .

In fact, the above definition of the geodesic subspaces (which is the one implicitly used in most of the works using PGA) may be too large and may not define a k -dimensional submanifold when there is a cut-locus. For instance, it is well known that geodesics of a flat square torus are either periodic or everywhere dense in a flat torus submanifold depending on whether the components of the initial velocity field have rational or irrational ratios. This means that the geodesic space generated by a single vector for which all ratio of coordinates are irrational (e.g. $w = (\pi, \pi^2, \dots, \pi^k)$) is filling the full k -dimensional flat torus. Thus all the 1-dimensional geodesic subspaces that have irrational ratio of all coordinates minimize the distance to any set of data points in a flat square torus of any dimension, which is not very interesting from the application point of view. In order to have a more meaningful definition and to guaranty the dimensionality of the geodesic subspace, we need to restrict the definition to the points of the geodesics that are minimizing the distance.

Definition 9 (Restricted Geodesic Submanifolds). *Let $x \in \mathcal{M}$ be a point of a Riemannian manifold and let $W_x = \{\sum_{i=1}^k \alpha_i w_i, \alpha \in \mathbb{R}^k\}$ be the k -dimensional linear subspace of $T_x \mathcal{M}$ generated a k -uplet $\{w_i\}_{1 \leq i \leq k} \in (T_x \mathcal{M})^k$ of tangent vectors at x .*

We call restricted geodesic submanifold $GS^(W_x)$ at x generated by the vector subspace W_x the submanifold of \mathcal{M} generated by the geodesics starting at x with tangent vectors $w \in W_x$, but up to the first cut-point of x only:*

$$GS^*(W_x) = GS^*(x, w_1, \dots, w_k) = \{\exp_x(w), w \in W_x \cap D(x)\},$$

where $D(x) \subset T_x \mathcal{M}$ is the maximal definition domain on which the exponential map is diffeomorphic.

It may not be immediately clear that the subspace we define that way is a submanifold of \mathcal{M} : since \exp_x is a diffeomorphism from $D(x) \subset T_x \mathcal{M}$ to $\mathcal{M}/\mathcal{C}(x)$ whose differential has full rank, its restriction to the open star-shape subset $W_x \cap D(x)$ of dimension k is a diffeomorphism from that subset to the restricted geodesic subspace $GS^*(W_x)$ which is thus an open submanifolds of dimension k of \mathcal{M} . However, this submanifold is generally not geodesically complete.

Theorem 6 (Restricted geodesic subspaces are limit of affine spans). *Points of the restricted geodesic submanifold $GS^*(W_x) = \{\exp_x(w), w \in W_x \cap D(x)\}$ are points of the affine span $\text{Aff}(x, x_1, \dots, x_k)$ parametrized by points at infinity of \mathcal{P}_k^* when the points $x_i = \exp_x(\eta w_i)$ are converging to x at first order along the tangent vectors w_i defining the k -dimensional subspace $W_x \subset T_x \mathcal{M}$.*

Proof. We first establish a useful formula exploiting the symmetry of the geodesics from x to $y \notin \mathcal{C}(x)$ with respect to time. Reverting time along a geodesic, we have: $\gamma_{(x, \overrightarrow{xy})}(t) = \gamma_{(y, \overrightarrow{yx})}(1-t)$, which means in particular that $\dot{\gamma}_{(x, \overrightarrow{xy})}(1) = -\dot{\gamma}_{(y, \overrightarrow{yx})}(0) = -\overrightarrow{yx}$. Since $\gamma_{(x, \overrightarrow{xy})}(t) = \exp_x(t\overrightarrow{xy})$, we obtain $\overrightarrow{yx} = -D \exp_x|_{\overrightarrow{xy}} \overrightarrow{xy}$. Now, we also have $(D \exp_x|_{\overrightarrow{xy}}) \cdot D \log_x|_y = \text{Id}$ because $\exp_x(\log_x(y)) = y$. Finally, $D \exp_x$ and $D \log_x$ have full rank on $\mathcal{M}/\mathcal{C}(x)$ since there is no conjugate point before the cut-locus, so that we can multiply by their inverse and we end up with:

$$\forall y \notin \mathcal{C}(x), \quad \overrightarrow{xy} = -D \log_x|_y \overrightarrow{yx}. \quad (8)$$

In order to work properly, let us first restrict to a convenient domain of \mathcal{M} : we consider a open geodesic ball $B(x_0, \epsilon)$ of radius ϵ centered at x_0 and we exclude all the points of \mathcal{M} which cut locus intersect this ball, or equivalently the cut-locus of all the points of this ball. We obtain an open domain $\mathcal{D}_\epsilon(x_0) = \mathcal{M}/\mathcal{C}(B(x_0, \epsilon))$ in which $\log_x(y)$ is well defined and smooth for all $x \in B(x_0, \epsilon)$ and all $y \in \mathcal{D}_\epsilon(x_0)$. Thanks to the symmetry of the cut-locus, $\log_y(x)$ is also well defined and smooth in the same conditions and Eq. (8) can be rephrased as:

$$\forall x \in B(x_0, \epsilon), y \in \mathcal{D}_\epsilon(x_0), \quad \overrightarrow{xy} = -D \log_x|_y \overrightarrow{yx}. \quad (9)$$

Let $\|w\|_\infty = \max_i \|w_i\|_{x_0}$ be the maximal length of the vectors w_i . For $\eta < \epsilon/\|w\|_\infty$, we have $\|\eta w_i\|_{x_0} \leq \eta \|w\|_\infty < \epsilon$, so that all the points $x_i = \exp_{x_0}(\eta w_i)$ belong to the open geodesic ball $B(x_0, \epsilon)$. Thus, $\log_x(x_i)$ and $\log_{x_i}(x)$ are well defined and smooth for any $x \in \mathcal{D}_\eta(x_0)$, and we can write the Taylor expansion in a normal coordinate system at x_0 :

$$\log_x(x_i) = \log_x(\exp_{x_0}(\eta w_i)) = \log_x(x_0) + \eta D \log_x|_{x_0} w_i + O(\eta^2).$$

Now let $x = \exp_{x_0}(w)$ (or $\log_{x_0}(x) = w$) with $w = \sum_{i=1}^k \alpha_i w_i \in \mathcal{D}_\epsilon(x_0)$. Using formula (9), we can write equivalently $\log_x(x_0) = -D \log_x|_{x_0} w$ and combined that with the above equation to obtain: $\log_x(x_i) = D \log_x|_{x_0}(\eta w_i - w) + O(\eta^2)$. Thus, the implicit equation $\mathfrak{M}_1(x, \lambda) = \sum_{i=0}^k \lambda_i \overrightarrow{xx_i} = 0$ is equivalent to $\sum_{i=0}^k \lambda_i (\eta w_i - w) = O(\eta^2)$, and showing that x is a point of the EBS amounts to find the normalized homogeneous coordinates λ satisfying

$$\sum_{i=0}^k (\lambda_i \eta - \alpha_i) w_i = O(\eta^2).$$

Taking $\lambda_i = \alpha_i$ for $1 \leq i \leq k$ and $\lambda_0 = \eta - (\sum_i \alpha_i)$ obviously satisfy this condition. With normalized coordinates, this writes: $\lambda_i = \alpha_i/\eta$ for $1 \leq i \leq k$

and $\lambda_0 = 1 - (\sum_i \alpha_i)/\eta$, so that we clearly see that it tends towards coordinates that sum up to zero (a point at infinity of \mathcal{P}_k^*). Thus, x is not a point of the EBS when η goes to zero, but since it is the limit of a series of points which are in it, it belongs to the closure of this set, which is the affine span. \square

We conjecture that the construction can be generalized using techniques from sub-Riemannian geometry to higher order derivatives when the first order derivative do not span a k -dimensional subspace. This would mean that we could also see some non-geodesic decomposition schemes as limit cases of barycentric subspaces, such as splines on manifolds that have been developed by Crouch and Leite [1995], Machado et al. [2010], Gay-Balmaz et al. [2012], Hinkle et al. [2014], Singh et al. [2015]. In fact, we will show in the next section shows that this is indeed the case for spheres where principal nested spheres developed by Jung et al. [2010, 2012] can be seen as a limit case of barycentric subspaces when they converge to a second-order jet.

4 Example on spheres

Intrinsic weighted averaging on spheres has been investigated in Buss and Fillmore [2001]. In particular, they have shown that for positive weights, there is a unique Fréchet mean if the points are within one hemisphere with at least one non-zero weight point not on the equator. In this section, we derive a similar result using different computations to exemplify how the barycentric subspaces are defined on spheres. We also provide show that the affine span is in generic conditions a great subsphere, and that it converges towards the Principal Nested Spheres when reference points converge toward a second-order jet.

4.1 Computing on spheres

We consider the unit sphere in dimension $n \geq 1$ embedded in \mathbb{R}^{n+1} so that points of $\mathcal{M} = \mathcal{S}_n$ are unit vectors of \mathbb{R}^{n+1} . The tangent space at x is the linear space of vectors orthogonal to x : $T_x \mathcal{S}_n = \{v \in \mathbb{R}^{n+1}, v^T x = 0\}$. The natural Riemannian metric on the unit sphere is inherited from the Euclidean metric of the embedding space. With these conventions, the Riemannian distance is the arc-length $d(x, y) = \arccos(x^T y) = \theta \in [0, \pi]$. Denoting $f(\theta) = 1/\text{sinc}\theta = \theta/\sin\theta$, the spherical exp and log maps are:

$$\exp_x(v) = \cos(\|v\|)x + \text{sinc}(\|v\|)v \quad (10)$$

$$\log_x(y) = f(\theta)(y - \cos\theta x) \quad \text{with } \theta = \arccos(x^T y). \quad (11)$$

Notice that $f(\theta)$ is a smooth function from $] -\pi; \pi[$ to \mathbb{R} that is always greater than one and is locally quadratic at zero: $f(\theta) = 1 + \theta^2/6 + O(\theta^4)$.

Using the orthogonal projection $v = (\text{Id} - xx^T)w$ of an unconstrained vector $w \in \mathbb{R}^{n+1}$ onto the tangent space $T_x \mathcal{S}_n$ we obtain a chart around a point $x \in \mathcal{S}_n$ where we can compute the gradient and Hessian of the squared-distance on the

sphere:

$$\nabla d_y^2(x) = -2f(\theta)(\text{Id} - xx^T)y = -2\log_x(y) \quad (12)$$

$$H_x(y) = 2vv^T + 2f(\theta)\cos\theta(\text{Id} - xx^T - vv^T) \quad (13)$$

The eigenvectors and eigenvalues of this matrix are easy to determine. By construction, x is an eigenvector with eigenvalue $\mu_0 = 0$. Then the vector v (or equivalently $\log_x(y) = \theta v$) is an eigenvector with eigenvalue $\mu_1 = 1$. To finish, every vector v which is orthogonal to these two vectors (i.e. orthogonal to the plane spanned by 0 , x and y) has eigenvalue $\mu_2 = f(\theta)\cos\theta = \theta \cot\theta$. This last eigenvalue is positive for $\theta \in [0, \pi/2[$, vanishes for $\theta = \pi/2$ and becomes negative for $\theta \in]\pi/2, \pi[$. We retrieve here the results of [Buss and Fillmore, 2001, lemma 2] expressed in a more general coordinate system.

4.2 $k + 1$ -pointed spheres

Let us now pick $k + 1$ points $X = \{x_0; \dots x_k\}$ on the sphere. We also denote by X the matrix of coordinates of the reference points $X = [x_0; \dots x_k]$. The cut locus of x_i is its antipodal point $-x_i$ so that the $(k + 1)$ -punctured manifold is $\mathcal{M}^*(x_0, \dots x_k) = \mathcal{S}_n / -X$. Denoting $\theta_i = \arccos(x_i^T x)$, we have $\log_x(x_i) = (\text{Id} - xx^T)f(\theta_i)x_i$, so that the first weighted moment is

$$\mathfrak{M}_1(x, \lambda) = (\text{Id} - xx^T) \sum_i \lambda_i f(\theta_i) x_i = (\text{Id} - xx^T) X F(X, x) \lambda$$

where $F(X, x) = \text{Diag}(f(\theta_i))$ is a diagonal matrix with entries that are always greater than one for $x \in \mathcal{M}^*(x_0, \dots x_k)$. Thus, we get the following expression for the reference matrix $Z(x)$, the covariance matrix of the reference points $\Sigma(x)$ and the and the Gram matrix $\Omega(x)$:

$$\begin{aligned} Z(x) &= (\text{Id} - xx^T) X F(X, x) \\ \Sigma(x) = Z(x) Z(x)^T &= (\text{Id} - xx^T) X F(X, x)^2 X^T (\text{Id} - xx^T) \\ \Omega(x) = Z(x)^T G(x) Z(x) &= F(X, x) X^T (\text{Id} - xx^T) X F(X, x). \end{aligned}$$

In all the above expressions, we recognize classical matrix equations, except for the scaling matrix $F(X, x)$ acting on homogeneous projective weights, which is non-stationary and non-linear in both X and x . In order to simplify all the computations, we introduce the change of coordinate system $\tilde{\lambda} = F(X, x)\lambda$ that we call renormalized weights. Since $F(X, x) = \text{Diag}(\theta_i / \sin\theta_i)$ is an invertible diagonal matrix, the original barycentric coordinates can be obtained by $\lambda = F(X, x)^{(-1)} \tilde{\lambda}$, or $\lambda_i = \tilde{\lambda}_i \sin\theta_i / \theta_i$.

Introducing the vector $\zeta(X, x) = F(X, x)^{(-1)} \mathbf{1} = [\sin\theta_0/\theta_0; \dots \sin\theta_k/\theta_k]^T$, the constraints $\mathbf{1}^T \lambda \neq 0$ on the original weights becomes $\zeta(X, x)^T \tilde{\lambda} \neq 0$ on the renormalized weights: it now excludes the hyperplane orthogonal to the vector $\zeta(X, x)$ (instead of $\mathbf{1}$) from the projective space \mathcal{P}_k . The reference and Gram matrices become $\tilde{Z}(x) = (\text{Id} - xx^T) X$ and $\tilde{\Omega}(x) = X^T (\text{Id} - xx^T) X$, which are now standard matrix expressions. In the following, we systematically work with the rescaled barycentric coordinates $\tilde{\lambda}$.

4.3 Exponential barycentric subspaces in generic conditions

The kernel of the matrix $\tilde{Z}(x)$ (or equivalently the Gram matrix $\tilde{\Omega}(x)$) determines the space of admissible (renormalized) weights $\tilde{\lambda}(x) = Ker(\tilde{Z}(x)) = Ker(\tilde{\Omega}(x))$. The solutions of the equation $\tilde{Z}(x)\tilde{\lambda} = 0$ under the constraint $\|x\| = 1$ are given by $(x^T X \tilde{\lambda})x = X \tilde{\lambda}$ or more explicitly $x = \pm X \tilde{\lambda} / \|X \tilde{\lambda}\|$. Thus, the point $x \in \mathcal{M}^*(X)$ has to belong to the Euclidean span of the reference vectors. Conversely, for any unit vector $x = X\alpha$ of the Euclidean span of X , we have $\tilde{Z}(x) = X - xx^T X = X - X\alpha\alpha^T X^T X$ so that $\tilde{Z}(x)\alpha = 0$ because $\|x\|^2 = \alpha^T X^T X \alpha = 1$. Thus $\tilde{\lambda} = \alpha$ are the renormalized barycentric coordinates of x (whenever x is not at the cut locus of the reference points so that $F(X, x)$ is invertible) and the non-normalized barycentric coordinates $\lambda = F(X, x)^{(-1)}\alpha$ are well defined. This shows that

$$EBS(X) = \text{Span}\{x_0, \dots, x_k\} \cap \mathcal{S}_n / \{-x_0, \dots, -x_k\}. \quad (14)$$

Notice that for each barycentric coordinates λ we have two antipodal solution points $x = \pm XF(X, x)\lambda$.

Using the renormalization principle, we can orthogonalize the reference points: let $X = USV^T$ be a singular value decomposition of the matrix of reference vectors. All the singular values s_i are positive since the reference vectors x_i are assumed to be linearly independent. Thus, the new renormalization $\check{\lambda} = SV^T \tilde{\lambda} = SV^T F(X, x)\lambda$ gives us the reference matrix $\check{Z}(x) = (\text{Id} - xx^T)U$. By definition of the singular value decomposition, the Euclidean span of X and U are the same, so that $EBS(U) = \text{Span}\{x_0, \dots, x_k\} \cap \mathcal{S}_n / -U$. This shows that the exponential barycentric subspace generated by the original points $X = [x_0; \dots, x_k]$ and the orthogonalized points $U = [u_0; \dots, u_k]$ are the same, except at the cut locus of all these points.

4.4 Affine spans in generic conditions

To obtain the affine span, we take the closure of the EBS, which incorporates the cut locus of the reference points:

$$\text{Aff}(X) = \text{Span}\{x_0, \dots, x_k\} \cap \mathcal{S}_n \quad (15)$$

Thus, for spherical data as for Euclidean data, the affine span only depend on the reference points through the point of the Grassmanian that they define. When the reference points are not linearly independent, the matrix X has one or more (say l) vanishing singular values. A singular value decomposition $X = USV^T$ shows that the value of $\tilde{\lambda}$ (and thus of $\lambda = F^{(-1)}(X, x)\tilde{\lambda}$) is in that case unconstrained in the vector space generated by the right singular vectors associated to the l vanishing singular values. Thus, the space of admissible weights at each point of the affine span is of dimension l , and the affine span itself is still the subsphere generated the Euclidean span of the reference vectors which is of dimension $k - l$.

Theorem 7 (Affine span in spheres). *The affine span $\text{Aff}(X)$ of $k+1$ different reference unit points $X = [x_0; \dots x_k]$ on the n -dimensional sphere \mathcal{S}_n provided with the canonical Euclidean metric of the embedding space \mathbb{R}^{n+1} is the great subsphere of dimension $\text{rank}(X) - 1$ that contains the reference points.*

4.5 Projection onto the affine span

A point unit vector of \mathbb{R}^{n+1} can be split in a unique way into its component within the Euclidean span of X and its Euclidean orthogonal complement. Among the vectors $\hat{x} = X\tilde{\lambda}$ of the Euclidean span of X , the closest one in the Euclidean sense is parametrized by the $\tilde{\lambda}$ solution of $X^T x = X^T X \tilde{\lambda}$, i.e. $\tilde{\lambda} = (X^T X)^{-1} X^T x$ when the Gram matrix $X^T X$ is full rank. When the Gram matrix is rank deficient, the smallest norm solution is given by the Moore-Penrose pseudo-inverse $\tilde{\lambda} = X^\dagger x$ and we can add any vector of the kernel of $X^T X$ without changing the point x . Thus, the component of x in $\text{Aff}_{\mathbb{R}^{n+1}}(X)$ is $\hat{x} = X X^\dagger x$ and the orthogonal component is $\check{x} = (\text{Id} - X X^\dagger)x$. Denoting $\phi = \arctan(\|\check{x}\|/\|\hat{x}\|) = \arccos(\|\hat{x}\|) = \arcsin(\|\check{x}\|)$, this amounts to the decomposition $x = \cos \phi x_{aff} + \sin \phi x^\perp$ on the sphere (if $\phi < \pi/2$) with

$$x_{aff} = \frac{\hat{x}}{\|\hat{x}\|} = \frac{X X^\dagger x}{\|X X^\dagger x\|} \quad \text{and} \quad x^\perp = \frac{\check{x}}{\|\check{x}\|} = \frac{(\text{Id} - X X^\dagger)x}{\|(\text{Id} - X X^\dagger)x\|}.$$

It turns out that the point x_{aff} is also the spherical projection of x (in the sense of the closest point) onto the spherical affine span of X (when $\hat{x} \neq 0$ so that it is defined).

4.6 Karcher barycentric subspaces

We turn in this section to the locus of local minima of the (normalized) weighted variance. Excluding the cut locus of the reference points from the analysis, we know that the critical points of the variance are the points of the EBS, so that the problem amounts to distinguish the minima from the maxima and saddle points (i.e non-degenerate and positive points according to definitions 8). For that, we compute the Hessian of the normalized variance. Using Eq.(25), we obtain :

$$H(x, \lambda) = \left(\sum_i \lambda_i \theta_i \cot \theta_i \right) (\text{Id} - x x^T) + \sum_i \lambda_i (1 - \theta_i \cot \theta_i) v_i v_i^T$$

As expected, x is an eigenvector with eigenvalue $\mu_0 = 0$ due to the projection on the tangent space at x . Any vector w of the tangent space at x (thus orthogonal to x) which is orthogonal to the affine span (and thus to the vectors v_i) is an eigenvector with eigenvalue $\mu_2(\lambda) = \left(\sum_i \lambda_i f(\theta_i) \cos \theta_i \right)$. Since the Euclidean affine span of X has $\text{rank}(X) \leq k+1$ dimensions, the multiplicity this eigenvalue is $n+1 - \text{rank}(X) \geq n-k$ if $x \in \text{Aff}(X)$ and $n - \text{rank}(X) \geq n-k-1$ otherwise. The last $\text{Rank}(X) - 1$ (resp $\text{Rank}(X)$) eigenvalues have associated eigenvectors within $\text{Aff}_{\mathbb{R}^{n+1}}(X)$.

Buss and Fillmore [2001] have shown that this Hessian matrix is positive definite for *positive weights* when the points are within one hemisphere with at least one non-zero weight point which is not on the equator. However, in our case, we are not interested by the positivity and definiteness of the Hessian $H(x, \lambda)$ for all the possible positive weights, but for the positive and negative weights which live in dual space of valid weights $\Lambda(x)$. This is actually a quite difficult algebraic geometry problem. However, simulation tests with random reference points X in general conditions show that $n - k$ non-zero eigenvalues of $H(x, \lambda(x))$ can be positive or negative at different points of the EBS. This simple test illustrate that the EBS on spheres is actually subdivided into different cells depending on the index of the critical point and that the positive points do not in general cover the full subsphere containing the reference points. Moreover, the frontiers of these cells do evolve when we move the reference points within the generated subsphere, contrarily to the affine span which consistently covers the whole subsphere. For subspace definition purposes, this suggests that the affine span might thus be the most interesting definition to work with.

4.7 Affine span with reference points coalescing at order 1

In the previous section, we assumed that all the reference points were distinct. We now investigate limit cases. We first assume that all the reference points coalesce to a single point $x_i = \exp_{x_0}(\epsilon w_i)$ along the tangent vectors w_i which are satisfying $x_0^T w_i = 0$ (to belong to the tangent space at x_0) and $\sum_i w_i = 0$. Denoting $X_0 = x_0 \mathbf{1}^T$, this amounts to say that we are following the curve $X_\epsilon = X_0 + \epsilon W$ in the space of affine spans, with $X_0^T W = 0$ and $W X_0^T = 0$. Solving the equation $\tilde{Z}(x) \tilde{\lambda} = 0$, we find that $x = \pm X_\epsilon \tilde{\lambda} / \alpha_\epsilon$ for some scalar factor α_ϵ that we can determine thanks to the orthogonality of X_0 and W : $\alpha_\epsilon = \|\mathbf{1}^T \tilde{\lambda} x_0 + \epsilon W \tilde{\lambda}\|$. Thus, we end-up with $x = x_0 + \epsilon W \tilde{\lambda} / \|\mathbf{1}^T \tilde{\lambda}\| + O(\epsilon^2)$, which shows that the space $EBS(X_0)$ is the intersection of the sphere with the Euclidean hyperplane going through x_0 generated by the vectors of W , minus the cut locus of x_0 . Thus, its completion $\text{Aff}(X)$ is once again the great subsphere generated by the completion of the geodesic subspace $GS(x_0, W)$.

4.8 Coalescence at order 2 and link with principal nested spheres

Principal nested spheres were proposed by Jung et al. [2010] and Jung et al. [2012] as a general framework for non-geodesic decomposition of high-dimensional spheres with applications to planar landmarks shape spaces. A subsphere \mathcal{A}_{n-1} of \mathcal{S}_n is defined as the set of points which are at a fixed distance $\theta \in (0, \pi/2]$ of a point $x \in \mathcal{S}_n$: $\mathcal{A}_{n-1}(x, \theta) = \{y \in \mathcal{S}_n / d(x, y) = \theta\}$. The subsphere $\mathcal{A}_{n-1}(x, \theta)$ can be viewed as the slice of \mathcal{S}_n by the n -dimensional affine hyperplane $P(x, \theta) = \{y \in \mathbb{R}^{n+1} / y^T x = \cos \theta\}$. Notice that the coordinates $(x, \cos \theta)$ of the affine hyperplane parametrize all the possible subspheres of dimension $n - 1$. In this process, the subsphere is not of radius one, unless one takes

$\theta = \pi/2$, in which case the hyperplane is passing through the origin. In this case, it is equivalent to the affine span, since any set of n non-collinear points x_1, \dots, x_n in the hyperplane $P(x, \theta)$ generate an affine span which is the great sphere $\mathcal{S}_{n-1} = \{y \in \mathbb{R}^{n+1} / \|x\| = 1, y^T x = 0\}$.

In order to figure out how smaller subspheres are related to affine spans, we consider the example of a circle of radius $\alpha \in]-1; 1[$ around the axis e_3 on the 3-sphere. It is described by the implicit equation $x^T e_3 = \sqrt{1 - \alpha^2}$ and the explicit equation: $x(\psi) = \alpha \cos(\psi)e_1 + \alpha \sin(\psi)e_2 + \sqrt{1 - \alpha^2}e_3$. Now, let us consider the barycentric subspace generated by the three points on that circle at angles $\psi_0 = 0, \psi_1 = \epsilon$ and $\psi_2 = -\epsilon$. The matrix of reference points is $X = [x(0), x(\epsilon), x(-\epsilon)]$. Using the change of coordinates $s = (\tilde{\lambda}_0 + \tilde{\lambda}_1 + \tilde{\lambda}_2)$, $u = (\tilde{\lambda}_1 - \tilde{\lambda}_2)\epsilon/s$ and $v = (\tilde{\lambda}_1 + \tilde{\lambda}_2)\epsilon^2/(2s)$, and setting the scaling factor $s = 1$ because we use homogeneous coordinates, we get that: $X\tilde{\lambda} = \alpha(1-v)e_1 + \alpha ue_2 + \sqrt{1 - \alpha^2}e_3 + O(\epsilon^3)$. Thus, the points of the affine span $x = X\tilde{\lambda}/\|X\tilde{\lambda}\|$ can only be in the hyperplane $x^T e_3 = \sqrt{1 - \alpha^2}$ when ϵ goes to zero, whose intersection with the sphere describes the original circle of radius α .

Iterating the process, one can generalize the above construction to subspheres of arbitrary dimensions. Thus, nested spheres can be seen as a limit of the affine span when the k reference points tend to a 2-jet. It would be interesting to determine which types of subspaces could be obtained by such limits for more general non-local and higher order jets.

5 Barycentric subspace analysis

This section generalizes principal component analysis itself. PCA can be viewed as the search for a sequence of nested linear spaces that best approximate the data at each level. In a Euclidean space, minimize the variance of the residues boils down to an independent optimization of orthogonal subspaces at each level of approximation, thanks to the Pythagorean theorem. This enables building each subspace of the sequence by adding (resp. subtracting) the optimal one-dimensional subspace iteratively in a forward (resp. backward) analysis. Of course, this property does not scale up to manifolds, for which the orthogonality of subspaces is not even well defined.

5.1 Flags of barycentric subspaces in manifolds

Damon and Marron [2013] have argued that the nestedness of approximation spaces is one of the most important characteristics for generalizing PCA to more general spaces. Barycentric subspaces can easily be nested, for instance by adding or removing one or several points at a time, to obtain a family of embedded submanifolds which generalizes flags of vector spaces.

A flag of a vector space V is a filtration of subspaces (an increasing sequence of subspaces, where each subspace is a proper subspace of the next): $\{0\} = V_0 \subset V_1 \subset V_2 \subset \dots \subset V_k = V$. Denoting $d_i = \dim(V_i)$ the dimension of the subspaces, we have $0 = d_0 < d_1 < d_2 < \dots < d_k = n$, where n is the dimension

of V . Hence, we must have $k \leq n$. A flag is *complete* if $d_i = i$, otherwise it is a *partial flag*. Notice that a linear subspace W of V can be identified to the partial flag $\{0\} \subset W \subset V$.

With barycentric subspaces of an n -dimensional manifold \mathcal{M} , an ordering of $n + 1$ distinct points $x_0 \prec x_1 \dots \prec x_n$ defines the filtration of subspaces: $BS(x_0) = \{x_0\} \subset \dots \subset BS(x_0, x_1, x_k) \subset \dots \subset BS(x_0, \dots, x_n)$. Notice that the 0-dimensional subspace is now a point in \mathcal{M} instead of the null vector in flags of vector spaces because we are in an affine setting. Grouping points together in the addition/removal process generates a partial flag of barycentric subspaces. Among the barycentric subspaces, the affine span seems to be the most interesting definition. Indeed, when the manifold $\mathcal{M}^*(x_0, \dots, x_k)$ is connected, the EBS of $n + 1$ distinct points covers the full manifold $\mathcal{M}^*(x_0, \dots, x_k)$, and its completion covers the original manifold: $\text{Aff}(x_0, \dots, x_n) = \mathcal{M}$. With the Fréchet or Karcher barycentric subspaces, we only generate a submanifold (the positive span) that does not cover the whole manifold in general, as we have seen with the example of spheres.

Definition 10 (Flags of affine spans in manifolds).

Let $x_0 \preceq x_1 \dots \preceq x_k$ be $k + 1 \leq n$ distinct and partially ordered points of \mathcal{M} . By partially ordered, we mean that two or more successive points can be considered as exchangeable ($x_i \sim x_{i+1}$). For a totally ordered set of points, we call flag of affine spans $FL(x_0 \prec x_1 \dots \prec x_k)$ the sequence of properly nested subspaces $FL_i(x_0 \prec x_1 \dots \prec x_k) = \text{Aff}(x_0, \dots, x_i)$ for $0 \leq i \leq k$. For partially ordered sets of points $x_0 \preceq x_1 \dots \preceq x_k$, subspaces in the sequence are only generated at strict ordering signs or at the end, so that all exchangeable points are always considered together.

A flag is said complete if it is totally ordered with $k = n$. We call pure subspace a flag of completely exchangeable points $FL(x_0 \sim x_1 \dots \sim x_k)$ because the sequence is reduced to the unique subspace $FL_k(x_0 \sim x_1 \dots \sim x_k) = \text{Aff}(x_0, \dots, x_k)$.

5.2 Forward and backward barycentric subspaces analysis

In Euclidean PCA, the flag of linear subspaces can be built in a forward way, by computing the best 0-th order approximation (the mean), then the best first order approximation (the first mode), etc. It can also be built backward, by removing the direction with the minimal residue from the current affine subspace. In a manifold, we can use similar forward and backward analysis, but they have no reason to give the same result.

With a forward analysis, we compute iteratively the flag of affine spans by adding one point at a time and keeping the previous ones fixed. Thus, we begin by computing the optimal barycentric subspace of dimension 0: $\text{Aff}(x_0) = \{x_0\}$. Since there is only one weight and it should be unit, the optimal point x_0 found by minimizing the unexplained variance is a Karcher mean. Adding a second point amounts to compute the geodesic passing through the mean that best approximate the data. Adding one more point now differ from PGA, unless

the three points coalesce to a single one. The procedure is continued point by point, which implies that the Fréchet mean always belong to the barycentric subspace. In practice, the forward analysis should be stopped when the variance of the residues reaches the noise level of the data, hopefully with k much lower than the dimension n of the embedding manifold, in order to realize an efficient dimension reduction.

The backward analysis consists in iteratively removing one dimension, thus one point in our case. One should theoretically start with a full set of points x_0, \dots, x_n which generates the full manifold and chose which one to remove. However, as all the sets of $n + 1$ distinct points generate the full manifold with the affine span, the optimization really begin with the set of n points x_0, \dots, x_{n-1} . Actually this should normally be the only time when we perform an optimization for the point positions, since one should afterward only test for which of the n points we should remove, and this optimization is particularly ill-posed and inefficient in large dimensional spaces!

In order to get around this problem, we may run a forward analysis until we reach the noise level of the data for a dimension $k \gg n$. Since the goal is only to characterize the optimal k -dimensional subspace, we may optimize the point positions at each step to better fit the data. Then, a backward sweep at the end only reorders the points if necessary by iteratively selecting the one that least increase the unexplained variance. With this process, there is no reason why the Fréchet mean should belong to the reference points (and even to any of the barycentric subspaces). For instance, if we have clusters of points, one expects the reference points to localize within these clusters rather than at the Fréchet mean.

5.3 Approximating data using a pure subspace

Let $Y = \{\hat{y}_i\}_{i=1}^N \in \mathcal{M}^N$ be N data points and $X = \{x_0, \dots, x_k\}$ be $k + 1$ distinct reference points. We assume in this analysis that each data point \hat{y}_i has almost surely one unique closest point $y_i(X)$ on the barycentric subspace. This is the situation that we observe for Euclidean spaces and for the sphere, and this should be the case for all the points outside the focal set of the barycentric subspace. This allows us to write the residual $r_i(X) = \text{dist}(\hat{y}_i, y_i(X))$ and to consider the minimization of the unexplained variance $\sigma_{out}^2(X) = \sum_j r_j^2(X)$. For a fixed number k of reference points $X = \{x_0, \dots, x_k\}$, this boils down to an optimization problem on \mathcal{M}^k , which can be achieved by standard techniques of optimization on manifolds (see e.g. Absil et al. [2008]). Here, it is not obvious that the canonical product Riemannian metric is the right metric to use, especially close to coincident points. In this case, one would like to consider switching to the space of (non-local) jets to guaranty the numerical stability of the solution. In practice, we may constraint on the distance between reference points to be larger than a threshold.

A second potential problem is the lack of identifiability: the minimum of the unexplained variance may not be unique. This is the case for instance in Euclidean spaces and on spheres for which every linearly independent k -uplet

of points in a given subspace parametrizes the same barycentric subspace. In a Euclidean space, this can be taken into account using a suitable polar or QR matrix factorization. In general manifolds, we expect that the absence of symmetries will break the multiplicity of this relationship (at least locally) thanks to the curvature. However, it can lead to very badly conditioned systems to solve from a numerical point of view for small curvatures.

A last problem is that the criterion we use here (the unexplained variance) is only valid for a pure subspace of fixed dimension, and considering a different dimension will lead in general to pure subspaces which cannot be described by a common subset of reference points. Thus, the forward and backward optimization of nested barycentric subspaces cannot lead to the simultaneous optimality of all the subspaces of a flag in general manifolds.

5.4 A criterion for hierarchies of subspaces: AUC on flags of affine spans

In order to obtain better properties, it is necessary to define a criterion which depends on the whole flag of subspaces and not on each of the subspaces independently. The simplest proposal for that is to sum the criterion of each of the subspaces for all the dimensions (accounting of course for the multiplicity in incomplete flags). In PCA, one often plots the unexplained variance as a function of the number of modes used to approximate the data. This curve should decrease as fast as possible from the variance of the data (for 0 modes) to 0 (for n modes). Summing the values at all steps amounts to compute the area under the curve, which is a standard way to quantify the decrease.

Given a totally ordered flag of affine subspaces $Fl(x_0 \prec x_1 \dots \prec x_k)$, we thus propose to optimize the AUC criterion:

$$AUC(Fl(x_0 \prec x_1 \dots \prec x_k)) = \sum_{i=0}^k \sigma_{out}^2(Fl_i(x_0 \prec x_1 \dots \prec x_k))$$

instead of the unexplained variance at order k . We could of course consider a complete flag but in practice it is often useful to stop at a dimension k much smaller than the possibly very high dimension n . The criterion is extended to more general partial flags by weighting the unexplained variance of each subspace by the number of (exchangeable) points that are added at each step. With this global criterion, the point x_i influences all the subspaces of the flag that are larger than $Fl_i(x_0 \prec x_1 \dots \prec x_k)$ but not the smaller subspaces. It turns out that optimizing this criterion results in the usual PCA up to mode k in a Euclidean space.

Theorem 8 (Euclidean PCA as an optimization in the flag space). *Let $\hat{Y} = \{\hat{y}_i\}_{i=1}^N$ be a set of N data points in \mathbb{R}^n . We denote as usual the mean by $\bar{y} = \frac{1}{N} \sum_{i=1}^N \hat{y}_i$ and the empirical covariance matrix by $\Sigma = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - \bar{y})(\hat{y}_i - \bar{y})^T$. Its spectral decomposition is denoted $\Sigma = \sum_{j=1}^n \sigma_j^2 u_j u_j^T$ with the eigenvalues*

sorted in decreasing order. We assume that the first $k + 1$ eigenvalues have multiplicity one, so that the order from σ_1 to σ_{k+1} is strict.

Then the partial flag of affine subspaces $Fl(x_0 \prec x_1 \dots \prec x_k)$ optimizing the AUC criterion:

$$AUC(Fl(x_0 \prec x_1 \dots \prec x_k)) = \sum_{i=0}^k \sigma_{out}^2(Fl_i(x_0 \prec x_1 \dots \prec x_k))$$

is totally ordered and can be parametrized by $x_0 = \bar{y}$, $x_i = x_0 + u_i$ for $1 \leq i \leq k$. The parametrization by points is not unique but the flag of subspaces which is generated is and is equal to the flag generated by the PCA modes up to mode k included.

The proof is not very difficult but relies on tedious computations which are summarized below. A more detailed proof is provided in the appendix B.

Proof. Given $k + 1$ reference points x_0, \dots, x_k in generic conditions, we first perform the QR decomposition $[x_1 - x_0; \dots; x_i - x_0; \dots; x_k - x_0] = QT$ using the Gram-Schmidt orthogonalization process, to obtain the $n \times k$ orthogonal matrix $Q = [q_1; \dots; q_k]$ and the $k \times k$ triangular superior matrix T with entries $t_{ij} = q_i^T(x_j - x_0)$ for $j \geq i$. The decomposition is unique when all the points are linearly independent. The affine span generated by the $k + 1$ reference points is thus $Aff(X = [x_0; x_1; \dots; x_k]) = \{x = x_0 + Q\alpha/\alpha \in \mathbb{R}^k\}$. The matrix T has no influence and may thus be freely chosen to be the identity, so that $Aff(X)$ is parametrized by x_0 and k orthonormal vectors $q_1 \dots q_k$. This partial orthonormal basis can be complemented by $n - k$ unit vector q_{k+1}, \dots, q_n to constitute a complete orthonormal basis of \mathbb{R}^n .

The key property of the QR decomposition is its stability under the removal of reference points: if we only consider the $i < k$ first reference points, then $Aff(X_i = [x_0; \dots; x_i])$ is parametrized by x_0 and the first i orthonormal vectors $Q_i = [q_0, \dots, q_i]$. Using these notations, the projection of a data point y on the subspace $Aff(X_i)$ is $Proj(y, Aff(X_i)) = x_0 + Q_i Q_i^T (y - x_0)$ and the unexplained variance (sum of squared residuals to $Aff(X_i)$) is: $\sigma_{out}^2(X_i) = \text{Tr}(W_i(\Sigma - (\bar{y} - x_0)(\bar{y} - x_0)^T))$, with $W_i = (\text{Id}_n - Q_i Q_i^T) = \sum_{j=i+1}^n q_j q_j^T$. Thus the AUC criterion can be written:

$$AUC(X_k) = \text{Tr}(\bar{W}(\Sigma - (\bar{y} - x_0)(\bar{y} - x_0)^T))$$

with $\bar{W} = \sum_{i=0}^k W_i = \sum_{i=1}^k i q_i q_i^T + (k + 1) \sum_{i=k+1}^n q_i q_i^T$. The minimum of this criterion over x_0 is obviously achieved for $x_0 = \bar{y}$. Optimizing one by one the q_i 's starting from q_1 and taking into account the orthonormality constraints, we find that each q_i (for $i \leq k$) should be an eigenvector of Σ . Assuming that all the eigenvalues σ_i^2 of Σ are different (so that they can be sorted in a strict order and the eigenvectors have multiplicity one), the optimal values are $q_i^* = u_{\pi(i)}$ for some permutation π of the number $1 \dots n$. The value of the AUC criterion at this critical point is $AUC([q_1^*, q_2^* \dots q_n^*]) = \sum_{i=q}^k i \sigma_{\pi(i)}^2 + (k + 1) \sum_{i=k+1}^n \sigma_{\pi(i)}^2$. Now, to finish, we can show that the permutation of two indices $\pi(i)$ and $\pi(j)$

give a lower (or equal) criterion when $\pi(i) < \pi(j)$. The global minimum is thus achieved for the identity permutation $\pi(i) = i$ for the indices $1 \leq i \leq k$. For the higher indices, any combination of the last $n - k$ eigenvectors of Σ gives the same value of the criterion. When some eigenvalues of Σ have a multiplicity larger than one, then the corresponding eigenvectors cannot be uniquely determined since they can be rotated within the eigenspace. We end-up in that case with a partial flag. \square

6 Discussion

We investigated in the paper several notions of subspaces in manifolds generalizing the notion of affine span in a Euclidean space. The Fréchet / Karcher / exponential barycentric subspaces are the nested locus of weighted Fréchet / Karcher / exponential barycenters with positive or negative weights summing up to 1. The affine span is the metric completion of the largest one (the EBS). It may be a multiply connected manifold with boundaries. The completeness of the affine span enables reconnecting part of the subspace that arrive from different directions at the cut-locus of reference points if needed. It also allows ensuring that there exists a closest point on the submanifold for data projection purposes, which is fundamental for dimension reduction purposes. The fact that modifying the power of the metric does not change the affine span is an unexpected stability result which suggests that the notion is quite central.

In the case of spheres, we have shown that the affine span encompasses both principal geodesic subspaces and principal nested subspheres as limit cases. It would be interesting to see if we can obtain other types of subspaces with higher order and non-local jets. The study of the subspaces that can be obtained with this technique should of course be pushed to other spaces where PCA is used. For instance, Eltzner et al. [2015] adaptively deforms the a product of spheres into a unique sphere to allow principal nested spheres (PNS) analysis. A quick look at the flat torus shows that the the cut-locus of $k + 1$ points in S_1^n divides the torus into k^n cells in which the affine span is a k -dimensional linear subspace. The subspaces generated in each cell are generally disconnected, but when points coalesce with each others into a jet, the number of cells decreases and at the limit we recover a single cell that contain a connected affine span. For a first order jet, we recover as expected the restricted geodesic subspace (here a linear subspace limited to the cut locus of the jet base-point), but higher order jets may generate more interesting curved subspaces that may better describe the data geometry.

The next practical step is obviously the implementation of generic algorithms to work with barycentric subspaces in general Riemannian manifolds. Example algorithms include: finding a point with given barycentric coordinates (there might be several so this has to be a local search); finding the closest point (and its coordinates) on the barycentric subspace; optimizing the reference points to minimize the residual error after projection of data points, etc. If such algorithms can be designed relatively simply for simple specific manifolds as we

have done here on the sphere, the generalization to general manifolds requires a study of the focal set of the barycentric subspaces. We conjecture that the focal set is a stratified set of zero measure in generic cases, but this remains to be established for guarantying the correct behavior of algorithms. Another difficulty will be non-identifiability of the subspace parameters for simple (constant curvature?) spaces like Euclidean and spheres, where the right parameter space is actually the k -Grassmanian. In more general manifolds, the curvature and the interaction with the cut-locus break the symmetry of the barycentric subspaces, but may lead to a poor numerical conditioning of the system and good renormalization techniques need to be designed to guaranty the numerical stability.

Finding the subspace that best explain the data can also be recast as a problem of optimization on manifolds. This raise the question of which metric should be considered on the space of barycentric subspaces. In this paper, we mainly see this space as the configuration space of $k + 1$ points in general position, with convergence to spaces of jets (including non-local jets) when several points coalesce. Such a construction was named Multispace by Olver [2001] in the context of symmetry-preserving numerical approximations to differential invariants. It is likely that similar techniques could be investigated to construct numerically stable implementations of barycentric subspaces of higher order parametrized by non-local jets, which are needed to optimize safely. Conversely, barycentric subspaces could help shedding a new light on the multispace construction for differential invariants.

Barycentric subspaces could probably be also used to extended methods like probabilistic PCA of Tipping and Bishop [1999] which was generalized to PGA by Zhang and Fletcher [2013]. A first easy step in that direction is to replace the reference points by reference distributions on the manifold and to look at the locus of weighted expected means. Interestingly, this procedure soften the constraints that we had in this paper about the cut locus. Thus, following Karcher [1977], reference distributions could be used in a mollifier smoothing approach to study the regularity of the barycentric subspaces.

For applications where data live on Lie groups, generalizing BS to more general non-Riemannian spaces like affine connection manifolds is a particularly appealing extension. In computational anatomy, for instance, deformations of shapes are lifted to a group of diffeomorphism for statistical purposes (see e.g. Lorenzi and Pennec [2013], Lorenzi et al. [2015]). All Lie groups can be provided with a bi-invariant symmetric Cartan-Schouten connection for which geodesics are the left and right translation of one-parameter subgroups. This provides the Lie group with an affine connection structure which may be metric or not. When the group is the direct product of compact and Abelian groups, it admits a bi-invariant metric for which the Cartan-Schouten connection is the natural Levy-Civita connection. Other groups do not admit any bi-invariant metric (this is the case for rigid transformations in more than 2 dimensions because of the semi-direct product), so that a Riemannian structure can only be left or right invariant but not both. However the bi-invariant Cartan-Schouten connection continues to exists, and one can design bi-invariant means using ex-

ponential barycenter as proposed by Pennec and Arsigny [2012]. Thus, we may still define exponential barycentric subspaces and affine spans in these affine connection spaces, the main difference being that the derivative of the log is not any more the Hessian of a distance function. This might considerably complexify the analysis of the generated subspaces.

The second topic of this paper concerns the generalization of PCA to manifolds using Barycentric Subspace Analysis (BSA). Damon and Marron [2013] argued that an interesting generalization of PCA should rely on nested sequence of relations, like embedded linear subspaces in the Euclidean space or embedded spheres in PNS. Barycentric subspaces can naturally be nested by adding or removing points or equivalently by setting the corresponding barycentric coordinate to zero. Thus we can easily generalize PCA to manifolds using a forward analysis by iteratively adding one or more points at a time. At the limit where points coalesce at the first order, this amounts to build a flag of (restricted) principal geodesic subspaces. Thus it generalizes the Principal Geodesic Analysis (PGA) of Fletcher et al. [2004], Sommer et al. [2013] when starting with a zeroth dimensional space (the Fréchet mean) and the Geodesic PCA (GPCA) of Huckemann and Ziezold [2006], Huckemann et al. [2010] when starting directly with a first order jet defining a geodesic. One can also design a backward analysis by starting with a large subspace and iteratively removing one or more points to define embedded subspaces. In the case of spheres, this corresponds to the Principal Nested Spheres procedure at the limit of points coalescing to a second order jet.

However, the greedy optimization of these forward/backward methods lead generally to different solutions which are not optimal for all subspace jointly. The key idea is to consider PCA as a joint optimization of the whole flag of subspaces instead of each subspace independently. In a Euclidean space, we showed that the sum of the unexplained variance with respect to all the subspaces of the hierarchy (the area under the curve of unexplained variance) is a proper criterion on the space of Euclidean flags. We proposed to extend this criterion to barycentric subspaces in manifolds, where an ordering of the reference points naturally defines a flag of nested barycentric subspaces. A similar idea could be used with other iterative least-squares methods like partial least-squares (PLS) which are also one-step at a time minimization methods.

Acknowledgements

This work was partially supported by the Erwin Schrödinger Institute in Vienna through a three-weeks stay in February 2015 during the program Infinite-Dimensional Riemannian Geometry with Applications to Image Matching and Shape. It was also partially supported by the Inria Associated team GeomStats between Asclepios and Holmes' lab at Stanford Statistics Dept, through a 3 month stay at Stanford in 2015. I would particularly like to thank Prof. Susan Holmes for fruitful discussions during the writing of the paper.

A Hessian of the Riemannian squared distance

This appendix details the notions of Riemannian geometry that are underlying the main paper. In particular, it investigates the Hessian of the Riemannian square distance whose definiteness controls the local regularity of the barycentric subspaces. This is exemplified on the Sphere.

A.1 Riemannian manifolds

A Riemannian manifold is a differential manifold provided with a smooth collection of scalar products $\langle \cdot | \cdot \rangle_x$ on each tangent space $T_x \mathcal{M}$ at point x of the manifold, called the Riemannian metric. In a chart, the metric is expressed by a symmetric positive definite matrix $G(x) = [g_{ij}(x)]$ where each element is given by the dot product of the tangent vector to the coordinate curves: $g_{ij}(x) = \langle \partial_i | \partial_j \rangle$. This matrix is called the *local representation of the Riemannian metric* in the chart x and the dot products of two vectors v and w in $T_x \mathcal{M}$ is now $\langle v | w \rangle_x = v^T G(x) w$.

A.1.1 Riemannian distance and geodesics

If we consider a curve $\gamma(t)$ on the manifold, we can compute at each point its instantaneous speed vector $\dot{\gamma}(t)$ (this operation only involves the differential structure) and its norm $\|\dot{\gamma}(t)\|_{\gamma(t)}$ to obtain the instantaneous speed (the Riemannian metric is needed for this operation). To compute the length of the curve, this value is integrated along the curve:

$$\mathcal{L}_a^b(\gamma) = \int_a^b \|\dot{\gamma}(t)\|_{\gamma(t)} dt = \int_a^b \left(\langle \dot{\gamma}(t) | \dot{\gamma}(t) \rangle_{\gamma(t)} \right)^{\frac{1}{2}} dt$$

The distance between two points of a connected Riemannian manifold is the minimum length among the curves joining these points. The curves realizing this minimum are called geodesics. Finding the curves realizing the minimum length is a difficult problem as any time-reparameterization is authorized. Thus one rather defines the metric geodesics as the critical points of the energy functional $\mathcal{E}(\gamma) = \frac{1}{2} \int_0^1 \|\partial_\gamma\|^2 dt$. It turns out that they also optimize the length functional but they are moreover parameterized proportionally to arc-length.

Let $[g^{ij}] = [g_{ij}]^{(-1)}$ be the inverse of the metric matrix (in a given coordinate system x) and $\Gamma_{jk}^i = \frac{1}{2} g^{im} (\partial_k g_{mj} + \partial_j g_{mk} - \partial_m g_{jk})$ the Christoffel symbols (using Einstein summation convention that implicit sum upon each index that appear up and down in the formula). The calculus of variations shows the geodesics are the curves satisfying the following second order differential system:

$$\ddot{\gamma}^i + \Gamma_{jk}^i \dot{\gamma}^j \dot{\gamma}^k = 0.$$

The fundamental theorem of Riemannian geometry states that on any Riemannian manifold there is a unique (torsion-free) connection which is compatible with the metric, called the Levi-Civita (or metric) connection. This connection

is determined in a local coordinate system through the Christoffel symbols: $\nabla_{\partial_i} \partial_j = \sum_k \Gamma_{ij}^k \cdot \partial_k$. For that choice of connection, shortest paths (geodesics) are auto-parallel curves ("straight lines").

In the following, we only consider the Levi-Civita connection and we assume that the manifold is geodesically complete, i.e. that the definition domain of all geodesics can be extended to \mathbb{R} . This means that the manifold has no boundary nor any singular point that we can reach in a finite time. As an important consequence, the Hopf-Rinow-De Rham theorem states that there always exists at least one minimizing geodesic between any two points of the manifold (i.e. whose length is the distance between the two points).

A.1.2 Normal coordinate systems

Let x be a point of the manifold that we consider as a local reference and v a vector of the tangent space $T_x \mathcal{M}$ at that point. From the theory of second order differential equations, we know that there exists one and only one geodesic $\gamma_{(x,v)}(t)$ starting from that point with this tangent vector. This allows to wrap the tangent space onto the manifold, or equivalently to develop the manifold in the tangent space along the geodesics (think of rolling a sphere along its tangent plane at a given point). The mapping $\exp_x(v) = \gamma_{(x,v)}(1)$ of each vector $v \in T_x \mathcal{M}$ to the point of the manifold that is reached after a unit time by the geodesic $\gamma_{(x,v)}(t)$ is called the *exponential map* at point x . Straight lines going through 0 in the tangent space are transformed into geodesics going through point x on the manifold and distances along these lines are conserved.

The exponential map is defined in the whole tangent space $T_x \mathcal{M}$ (since the manifold is geodesically complete) but it is generally one-to-one only locally around 0 in the tangent space (i.e. around x in the manifold). In the sequel, we denote by $\vec{x}\vec{y} = \log_x(y)$ the inverse of the exponential map: this is the smallest vector (in norm) such that $y = \exp_x(\vec{x}\vec{y})$. It is natural to search for the maximal domain where the exponential map is a diffeomorphism. If we follow a geodesic $\gamma_{(x,v)}(t) = \exp_x(t v)$ from $t = 0$ to infinity, it is either always minimizing all along or it is minimizing up to a time $t_0 < \infty$ and not any more after (thanks to the geodesic completeness). In this last case, the point $\gamma_{(x,v)}(t_0)$ is called a *cut point* and the corresponding tangent vector $t_0 v$ a *tangential cut point*. The set of tangential cut points at x is called the *tangential cut locus* $C(x) \in T_x \mathcal{M}$, and the set of cut points of the geodesics starting from x is the *cut locus* $\mathcal{C}(x) = \exp_x(C(x)) \in \mathcal{M}$. This is the closure of the set of points where several minimizing geodesics starting from x meet. On the sphere $\mathcal{S}_2(1)$ for instance, the cut locus of a point x is its antipodal point and the tangential cut locus is the circle of radius π .

The maximal bijective domain of the exponential chart is the domain $D(x)$ containing 0 and delimited by the tangential cut locus ($\partial D(x) = C(x)$). This domain is connected and star-shaped with respect to the origin of $T_x \mathcal{M}$. Its image by the exponential map covers all the manifold except the cut locus, which has a null measure. Moreover, the segment $[0, \vec{x}\vec{y}]$ is mapped to the unique minimizing geodesic from x to y : geodesics starting from x are straight

lines, and the distance from the reference point are conserved. This chart is somehow the “most linear” chart of the manifold with respect to the reference point x .

When the tangent space is provided with an orthonormal basis, this is called *an normal coordinate systems at x* . A set of normal coordinate systems at each point of the manifold realize an atlas which allows to work very easily on the manifold. The implementation of the exponential and logarithmic maps (from now on \exp and \log) is indeed the basis of programming on Riemannian manifolds, and we can express using them practically all the geometric operations needed for statistics [Pennec, 2006] or image processing [Pennec et al., 2006].

The size of the maximal definition domain is quantified by the *injectivity radius* $\text{inj}(\mathcal{M}, x) = \text{dist}(x, \mathcal{C}(x))$, which is the maximal radius of centered balls in $T_x\mathcal{M}$ on which the exponential map is one-to-one. The injectivity radius of the manifold $\text{inj}(\mathcal{M})$ is the infimum of the injectivity over the manifold. It may be zero, in which case the manifold somehow tends towards a singularity (think e.g. to the surface $z = 1/\sqrt{x^2 + y^2}$ as a sub-manifold of \mathbb{R}^3).

In a Euclidean space, normal coordinate systems are realized by orthonormal coordinates system translated at each point: we have in this case $\vec{x}\vec{y} = \log_x(y) = y - x$ and $\exp_x(\vec{v}) = x + \vec{v}$. This example is more than a simple coincidence. In fact, most of the usual operations using additions and subtractions may be reinterpreted in a Riemannian framework using the notion of *bipoint*, an antecedent of vector introduced during the 19th Century. Indeed, vectors are defined as equivalent classes of bipoints in a Euclidean space. This is possible because we have a canonical way (the translation) to compare what happens at two different points. In a Riemannian manifold, we can still compare things locally (by parallel transportation), but not any more globally. This means that each “vector” has to remember at which point of the manifold it is attached, which comes back to a bipoint.

A.2 Hessian of the squared distance

A.2.1 Computing the differential of the Riemannian log

On $\mathcal{M}/\mathcal{C}(y)$, the gradient of the squared distance $d_y^2(x) = \text{dist}^2(x, y)$ with respect to the fixed point y is well defined and is equal to $\nabla d_y^2(x) = -2 \log_x(x_i)$. The Hessian operator $\nabla^2 f(x)$ from $T_x\mathcal{M}$ to $T_x\mathcal{M}$ is the covariant derivative of the gradient, defined by the identity $\nabla^2 f(v) = \nabla_v(\nabla f)$. In a chart (for instance a normal coordinate system at point x), the Hessian operator of the squared distance is thus

$$\nabla^2 d_y^2(x) = -2(D_x \log_x(x_i))$$

The points x and $y = \exp_x(v)$ are called conjugate if $D \exp_x(v)$ is singular. It is known that the cut point (if it exists) occurs at or before the first conjugate point along any geodesic [Lee, 1997]. Thus, $D \exp_x(v)$ has full rank inside the tangential cut-locus of x . This is in essence why there is a well posed inverse function $\vec{x}\vec{y} = \log_x(y)$, called the Riemannian log, which is continuous and differentiable everywhere except at the cut locus of x . Moreover, its differential

can be computed easily: since $\exp_x(\log_x(y)) = y$, we have $D \exp_x|_{\overrightarrow{xy}} D \log_x(y) = \text{Id}$, so that

$$D \log_x(y) = \left(D \exp_x|_{\overrightarrow{xy}} \right)^{-1} \quad (16)$$

is well defined and of full rank on $\mathcal{M}/C(x)$.

We can also see the Riemannian log $\log_x(y) = \overrightarrow{xy}$ as a function of the foot-point x , and differentiating $\exp_x(\log_x(y)) = y$ with respect to it gives: $D_x \exp_x|_{\overrightarrow{xy}} + D \exp_x|_{\overrightarrow{xy}} \cdot D_x \log_x(y) = 0$. Once again, we obtain a well defined and full rank differential for $x \in \mathcal{M}/C(y)$:

$$D_x \log_x(y) = - \left(D \exp_x|_{\overrightarrow{xy}} \right)^{-1} D_x \exp_x|_{\overrightarrow{xy}}. \quad (17)$$

The Hessian of the squared distance can thus be written:

$$\frac{1}{2} \nabla^2 d_y^2(x) = -D_x \log_x(x_i) = \left(D \exp_x|_{\overrightarrow{xy}} \right)^{-1} D_x \exp_x|_{\overrightarrow{xy}}.$$

If we notice that $J_0(t) = D \exp_x|_{t\overrightarrow{xy}}$ (respectively $J_1(t) = D_x \exp_x|_{t\overrightarrow{xy}}$) are actually matrix Jacobi field solutions of the Jacobi equation $\ddot{J}(t) + R(t)J(t) = 0$ with $J_0(0) = 0$ and $\dot{J}_0(0) = \text{Id}_n$ (respectively $J_1(0) = \text{Id}_n$ and $\dot{J}_1(0) = 0$), we see that the above formulation of the Hessian operator is equivalent to the one of Villani [2011][Equation 4.2]: $\frac{1}{2} \nabla^2 d_y^2(x) = J_0(1)^{(-1)} J_1(1)$.

A.2.2 Taylor expansion of the Riemannian log

In order to better figure out what the dependence of the Hessian of the squared Riemannian distance on curvature, we compute here the Taylor expansion of the Riemannian log function. Following Brewin [2009], we consider a normal coordinate system centered at x and $x_v = \exp_x(v)$ a variation of the point x . We denote by $R_{ihjk}(x)$ the coefficients of the curvature tensor at x and by ϵ a conformal gauge scale that encodes the size of the path in terms of $\|v\|_x$ and $\|\overrightarrow{xy}\|_x$ normalized by the curvature (see Brewin [2009] for details).

In a normal coordinate system centered at x , we have the following Taylor expansion of the metric tensor coefficients:

$$\begin{aligned} g_{ab}(v) = & g_{ab} - \frac{1}{3} R_{cabd} v^c v^d - \frac{1}{6} \nabla_e R_{cabd} v^e v^c v^d \\ & + \left(-\frac{1}{20} \nabla_e \nabla_f R_{cabd} + \frac{2}{45} R_{cad}^g R_{ebf}^h \delta_{gh} \right) v^c v^d v^e v^f + O(\epsilon^5). \end{aligned} \quad (18)$$

A geodesic joining point z to point $z + \delta z$ has tangent vector:

$$\begin{aligned} [\log_z(z + \Delta z)]^a = & \Delta z^a + \frac{1}{3} z^b \Delta z^c \Delta z^d R_{cbd}^a + \frac{1}{12} z^b z^c \Delta z^d \Delta z^e \nabla_d R_{bce}^a \\ & + \frac{1}{6} z^b z^c \Delta z^d \Delta z^e \nabla_b R_{dce}^a + \frac{1}{24} z^b z^c \Delta z^d \Delta z^e \nabla^a R_{bdce} \\ & + \frac{1}{12} z^b \Delta z^c \Delta z^d \Delta z^e \nabla_c R_{dbe}^a + O(\epsilon^4). \end{aligned}$$

Using $z = v$ and $z + \Delta z = \vec{x}\vec{y}$ (i.e. $\Delta z = \vec{x}\vec{y} - v$) in a normal coordinate system centered at x , and keeping only the first order terms in v , we obtain the first terms of the series development of the log:

$$[\log_{x+v}(y)]^a = \vec{x}\vec{y}^a - v^a + \frac{1}{3}R_{cbd}^a v^b \vec{x}\vec{y}^c \vec{x}\vec{y}^d + \frac{1}{12}\nabla_c R_{dbe}^a v^b \vec{x}\vec{y}^c \vec{x}\vec{y}^d \vec{x}\vec{y}^e + O(\epsilon^4). \quad (19)$$

Thus, the differential of the log with respect to the foot point is:

$$- [D_x \log_x(y)]_b^a = \delta_b^a - \frac{1}{3}R_{cbd}^a \vec{x}\vec{y}^c \vec{x}\vec{y}^d - \frac{1}{12}\nabla_c R_{dbe}^a \vec{x}\vec{y}^c \vec{x}\vec{y}^d \vec{x}\vec{y}^e + O(\epsilon^3) \quad (20)$$

Since we are in a normal coordinate system, the zeroth order term is the identity matrix, like in the Euclidean space, and the first order term vanishes. The Riemannian curvature tensor appear in the second order term and its covariant derivative in the third order term. The important point here is to see that the curvature is the leading term that makes this matrix departing from the identity (i.e. the Euclidean case) and which may lead to the non invertibility of the differential.

A.3 Example on spheres

We consider the unit sphere in dimension $n \geq 1$ embedded in \mathbb{R}^{n+1} and we represent points of $\mathcal{M} = \mathcal{S}_n$ as unit vectors in \mathbb{R}^{n+1} . The tangent space at x is naturally represented by the linear space of vectors orthogonal to x : $T_x \mathcal{S}_n = \{v \in \mathbb{R}^{n+1}, v^T x = 0\}$. The natural Riemannian metric on the unit sphere is inherited from the Euclidean metric of the embedding space \mathbb{R}^{n+1} . With these conventions, the Riemannian distance is the arc-length $d(x, y) = \arccos(x^T y) = \theta \in [0, \pi]$. Denoting $f(\theta) = 1/\text{sinc}(\theta) = \theta/\sin(\theta)$, the spherical exp and log maps are:

$$\exp_x(v) = \cos(\|v\|)x + \text{sinc}(\|v\|)v \quad (21)$$

$$\log_x(y) = f(\theta)(y - \cos(\theta)x) \quad \text{with} \quad \theta = \arccos(x^T y). \quad (22)$$

Notice that $f(\theta)$ is a smooth function from $] -\pi; \pi[$ to \mathbb{R} that is always greater than one and is locally quadratic at zero: $f(\theta) = 1 + \theta^2/6 + O(\theta^4)$.

A.3.1 Hessian of the squared distance on the sphere

To compute the gradient and Hessian of functions on the sphere, we first need a chart in a neighborhood of a point $x \in \mathcal{S}_n$. We consider the unit vector $x_v = \exp_x(v)$ which is a variation of x parametrized by the tangent vector $v \in T_x \mathcal{S}_n$ (i.e. verifying $x^T v = 0$). In order to extend this mapping to the embedding space to simplify computations, we consider that v is the orthogonal projection of an unconstrained vector $w \in \mathbb{R}^{n+1}$ onto the tangent space at x : $v = (\text{Id} - xx^T)w$. Using the above formula for the exponential map, we get at first order $x_v = x - v + O(\|v\|^2)$ in the tangent space or $x_v = x + (\text{Id} - xx^T)w + O(\|w\|^2)$ in the embedding space.

It is worth verifying first that the gradient of the squared distance $\theta^2 = d_y^2(x) = \arccos^2(x^\top y)$ is indeed $\nabla d_y^2(x) = -2\log_x(y)$. We consider the variation $x_w = \exp_x((\text{Id} - xx^\top)w) = x + (\text{Id} - xx^\top)w + O(\|w\|^2)$. Because $D_x \arccos(y^\top x) = -y^\top / \sqrt{1 - (y^\top x)^2}$, we get:

$$D_w \arccos^2(x_w^\top y) = \frac{-2\theta}{\sin \theta} y^\top (\text{Id} - xx^\top) = -2f(\theta) y^\top (\text{Id} - xx^\top),$$

and the gradient is as expected:

$$\nabla d_y^2(x) = -2f(\theta)(\text{Id} - xx^\top)y = -2\log_x(y). \quad (23)$$

To obtain the Hessian, we now compute the Taylor expansion of $\log_{x_w}(y)$. First, we have

$$f(\theta_w) = f(\theta) - \frac{f'(\theta)}{\sin \theta} y^\top (\text{Id} - xx^\top)w + O(\|w\|^2),$$

with $f'(\theta) = (1 - f(\theta) \cos \theta) / \sin \theta$. Thus, the first order Taylor expansion of $\log_{x_w}(y)$ is:

$$\log_{x_w}(y) = \left(f(\theta) - \frac{f'(\theta)}{\sin \theta} y^\top (\text{Id} - xx^\top)w \right) (\text{Id} - xx^\top - (\text{Id} - xx^\top)wx^\top - xw^\top (\text{Id} - xx^\top))y + O(\|w\|^2)$$

so that

$$-2D_w \log_{x_w}(y) = \frac{f'(\theta)}{\sin \theta} (\text{Id} - xx^\top)yy^\top (\text{Id} - xx^\top) - f(\theta) (x^\top y \text{Id} + xy^\top) (\text{Id} - xx^\top)$$

Now, since we have computed the derivative in the embedding space, we have obtained the Hessian with respect to the flat connection of the embedding space, which exhibits a non-zero normal component. In order to obtain the Hessian with respect to the connection of the sphere, we need to project back on $T_x \mathcal{S}_n$ (i.e. multiply by $(\text{Id} - xx^\top)$ on the left) and we obtain:

$$\begin{aligned} \frac{1}{2} H_x(y) &= \left(\frac{1 - f(\theta) \cos \theta}{\sin^2 \theta} \right) (\text{Id} - xx^\top) yy^\top (\text{Id} - xx^\top) + f(\theta) \cos \theta (\text{Id} - xx^\top) \\ &= (\text{Id} - xx^\top) \left((1 - f(\theta) \cos \theta) \frac{yy^\top}{\sin^2 \theta} + f(\theta) \cos \theta \text{Id} \right) (\text{Id} - xx^\top), \end{aligned}$$

To simplify this expression, we note that $\|(\text{Id} - xx^\top)y\|^2 = \sin \theta$, so that $v = \frac{(\text{Id} - xx^\top)y}{\sin \theta} = \frac{\log_x(y)}{\theta}$ is a unit vector of the tangent space at x (for $y \neq y$ so that $\theta > 0$). Using this unit vector and the intrinsic parameters $\log_x(y)$ and $\theta = \|\log_x(y)\|$, we can rewrite the Hessian:

$$\frac{1}{2} H_x(y) = f(\theta) \cos \theta (\text{Id} - xx^\top) + \left(\frac{1 - f(\theta) \cos \theta}{\theta^2} \right) \log_x(y) \log_x(y)^\top \quad (24)$$

$$= vv^\top + f(\theta) \cos \theta (\text{Id} - xx^\top - vv^\top) \quad (25)$$

The eigenvectors and eigenvalues of this matrix are now very easy to determine. By construction, x is an eigenvector with eigenvalue $\mu_0 = 0$. Then the vector v (or equivalently $\log_x(y) = f(\theta)(\text{Id} - xx^T)y = \theta v$) is an eigenvector with eigenvalue $\mu_1 = 1$. Lastly, every vector v which is orthogonal to these two vectors (i.e. orthogonal to the plane spanned by $0, x$ and y) has eigenvalue $\mu_2 = f(\theta) \cos \theta = \theta \cot \theta$. This last eigenvalue is positive for $\theta \in [0, \pi/2[$, vanishes for $\theta = \pi/2$ and becomes negative for $\theta \in]\pi/2, \pi[$. We retrieve here the results of [Buss and Fillmore, 2001, lemma 2] expressed in a more general coordinate system.

B PCA as an optimization on the flag manifold

This appendix details in length the proof that the flag of linear subspaces found by PCA optimizes the Area-Under-the-Curve (AUC) criterion in a Euclidean space.

B.1 A QR decomposition of the reference matrix

Let $X = [x_0; \dots x_k]$ be a matrix of $k + 1$ independent reference points in \mathbb{R}^n . Following the notations of the main paper, we write the reference matrix

$$Z(x) = [x - x_0; \dots x - x_k] = x \mathbf{1}_{k+1}^T - X.$$

The affine span $\text{Aff}(X)$ is the locus of points x satisfying $Z(x)\lambda = 0$ i.e. $x = X\lambda / (\mathbf{1}_{k+1}^T \lambda)$. Here, working with the barycentric weights is not so convenient, and in view of the principal component analysis, we prefer to work with a variant of the QR decomposition using the Gram-Schmidt orthogonalization process.

Choosing x_0 as the pivot point, we iteratively decompose $X - x_0 \mathbf{1}_{k+1}^T$ to find an orthonormal basis of the affine span of X . For convenience, we define the zeroth vectors $v_0 = q_0 = 0$. The first axis is defined by $v_1 = x_1 - x_0$, or by the unit vector $q_1 = v_1 / \|v_1\|$. Next, we project the second direction $x_2 - x_0$ onto $\text{Aff}(x_0, x_1) = \text{Aff}(x_0, x_0 + e_1)$: the orthogonal component $v_2 = (\text{Id} - e_1 e_1^T)(x_2 - x_0)$ is described by the unit vector $q_2 = v_2 / \|v_2\|$. The general iteration is then (for $i \geq 1$):

$$v_i = \left(\text{Id} - \sum_{j=0}^{i-1} e_j e_j^T \right) (x_i - x_0), \quad \text{and} \quad q_i = v_i / \|v_i\|.$$

Thus, we obtain the decomposition:

$$\begin{aligned} X &= x_0 \mathbf{1}_{k+1}^T + QT \\ Q &= [q_0; q_1; \dots q_k] \\ T &= \begin{bmatrix} q_0^T(x_0 - x_0) & q_0^T(x_1 - x_0) & q_0^T(x_2 - x_0) & \dots & q_0^T(x_k - x_0) \\ 0 & q_1^T(x_1 - x_0) & q_1^T(x_2 - x_0) & \dots & q_1^T(x_k - x_0) \\ 0 & 0 & q_2^T(x_2 - x_0) & \dots & q_2^T(x_k - x_0) \\ 0 & 0 & \dots & \dots & \dots \\ 0 & 0 & \dots & \dots & q_k^T(x_k - x_0) \end{bmatrix} \end{aligned}$$

With this affine variant of the QR decomposition, the $(k+1) \times (k+1)$ matrix T is triangular superior with vanishing first row and first column (since $q_0 = 0$). The $n \times (k+1)$ matrix Q also has a first null vector before the usual k orthonormal vectors in its $k+1$ columns. The decomposition into matrices of this form is unique when we assume that all the points x_0, \dots, x_k are linearly independent. This means that we can parametrize the matrix X by the orthogonal (aside the first vanishing column) matrix Q and the triangular (with first row and column zero matrix) T .

In view of PCA, it is important to notice that the decomposition is stable under the addition/removal of reference points. Let $X_i = [x_0; \dots; x_i]$ be the matrix of the first $i+1$ reference points (we assume $i < k$ to simplify here) and $X_i = x_0 \mathbf{1}_{i+1}^T + Q_i T_i$ his QR factorization. Then, the matrix Q_i is made of the first $i+1$ columns of Q and the matrix T_i is the upper $(i+1) \times (i+1)$ bloc of the upper triangular matrices T .

B.2 Optimizing the k -dimensional subspace

With our decomposition, we can now write any point of $x \in \text{Aff}(X)$ as the base-point x_0 plus any linear combination of the vectors q_i : $x = x_0 + Q\alpha$ with $\alpha \in \mathbb{R}^{k+1}$. The projection of a point y on $\text{Aff}(X)$ is thus parametrized by the $k+1$ dimensional vector α that minimizes the (squared) distance $d(x, y)^2 = \|x_0 + Q\alpha - y\|^2$. Notice that we have $Q^T Q = \text{Id}_{k+1} - e_1 e_1^T$ (here e_1 is the first basis vector of the embedding space \mathbb{R}^{k+1}) so that $Q^\dagger = Q^T$. The null gradient of this criterion implies that α is solving $Q^T Q \alpha = Q^T (y - x_0)$, i.e. $\alpha = Q^\dagger (y - x_0) = Q^T (y - x_0)$. Thus, the projection of y on $\text{Aff}(X)$ is

$$\text{Proj}(y, \text{Aff}(X)) = x_0 + Q Q^T (y - x_0),$$

and the residue is

$$r^2(y) = \|(\text{Id}_n - Q Q^T)(y - x_0)\|^2 = \text{Tr}((\text{Id}_n - Q Q^T)(y - x_0)(y - x_0)^T)$$

Accounting now for the N data points $Y = \{y_i\}_{i=1}^N$, and denoting as usual $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$ and $\Sigma = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})(y_i - \bar{y})^T$, the unexplained variance is:

$$\sigma_{out}^2(X) = \text{Tr}((\text{Id}_n - Q Q^T)(\Sigma - (\bar{y} - x_0)(\bar{y} - x_0)^T)).$$

In this formula, we see that the value of the upper triangular matrix T does not appear and can thus be chosen freely. The point x_0 that minimizes the unexplained variance is evidently $x_0 = \bar{y}$. To determine the matrix Q , we diagonalize the empirical covariance matrix to obtain the spectral decomposition $\Sigma = \sum_{j=1}^n \sigma_j^2 u_j u_j^T$ where by convention, the eigenvalues are sorted in decreasing order. The remaining unexplained variance $\sigma_{out}^2(X) = \text{Tr}((\text{Id}_n - (U^T Q)(U^T Q)^T) \text{Diag}(\sigma_i^2))$ reaches its minimal value $\sum_{i=k+1}^n \sigma_i^2$ for $[q_1, \dots, q_k] = [u_1, \dots, u_k] R$ where R is any $k \times k$ orthogonal matrix. Here, we see that the solution is unique in terms of subspaces (we have $\text{Span}(q_1, \dots, q_k) = \text{Span}(u_1, \dots, u_k)$ whatever orthogonal matrix R we choose) but not in terms of the matrix Q . In particular, the matrix $X = [\bar{y}, \bar{y} + u_1, \dots, \bar{y} + u_k]$ is one of the matrices describing the optimal subspace but the order of the vectors is not prescribed.

B.3 The AUC criterion

In PCA, one often plots the unexplained variance as a function of the number of modes used to approximate the data. This curve should decrease as fast as possible from the variance of the data (for 0 modes) to 0 (for n modes). Summing the values at all steps amount to compute the area under the curve, which is a standard way to quantify the decrease. We show in this section that the optimal flag of subspaces (up to dimension k) that optimize the AUC criterion is precisely the result of the PCA analysis.

As previously, we consider $k + 1$ points x_i but they are now ordered. We denote by $X_i = [x_0; \dots x_i]$ the matrix of the first $i+1$ columns of $X = [x_0; \dots x_k]$. The flag generated by X is thus

$$Aff(X_0) = \{x_0\} \subset Aff(X_1) \subset \dots \subset Aff(X) \subset \mathbb{R}^n.$$

The QR decomposition of X gives k orthonormal unit vectors $q_1 \dots q_k$ which can be complemented by $n - k$ unit vector $q_{k+1}, \dots q_n$ to constitute an orthonormal basis of \mathbb{R}^n . Using this extended basis, we can write:

$$\sigma_{out}^2(X) = \text{Tr}(W(\Sigma - (\bar{y} - x_0)(\bar{y} - x_0)^T))$$

with $W = (\text{Id}_n - QQ^T) = \sum_{j=k+1}^n q_j q_j^T$. Since the decomposition is stable under the removal of reference points, we have the QR factorization $X_i = x_0 \mathbb{1}_{i+1}^T + Q_i T_i$ with $Q_i = [q_0; \dots q_i]$ and we can write the unexplained variance for the subspace $Aff(X_i)$ as:

$$\sigma_{out}^2(X_i) = \text{Tr}(W_i(\Sigma - (\bar{y} - x_0)(\bar{y} - x_0)^T))$$

with $W_i = (\text{Id}_n - Q_i Q_i^T) = \sum_{j=i+1}^n q_j q_j^T$. Plugging this value into the criterion $AUC(X) = \sum_{i=0}^k \sigma_{out}^2(X_i)$, we get:

$$AUC(X_k) = \text{Tr}(\bar{W}(\Sigma - (\bar{y} - x_0)(\bar{y} - x_0)^T))$$

with

$$\bar{W} = \sum_{i=0}^k W_i = \sum_{i=0}^k (\text{Id}_n - Q_i Q_i^T) = \sum_{i=0}^k \sum_{j=i+1}^n q_j q_j^T = \sum_{i=1}^k i q_i q_i^T + (k+1) \sum_{i=k+1}^n q_i q_i^T.$$

B.4 PCA optimizes the AUC criterion

The minimum over x_0 is achieved as before for $x_0 = \bar{y}$ and the AUC for this value is now parametrized only by the matrix Q :

$$AUC(Q) = \text{Tr}(U^T W_k U \text{Diag}(\sigma_i^2)) = \sum_{i=1}^k i q_i^T \Sigma q_i + (k+1) \sum_{i=k+1}^n q_i^T \Sigma q_i.$$

Assuming that the the first $k + 1$ eigenvalues σ_i^2 ($1 \leq i \leq k + 1$) of Σ are all different (so that they can be sorted in a strict order), we claim that the

optimal unit orthogonal vectors are $q_i = u_i$ for $1 \leq i \leq k$ and $[q_{k+1}, \dots, q_n] = [u_{k+1}, \dots, u_n]R$ where $R \in O(n-k)$ is any orthogonal matrix.

In order to simplify the proof, we start by assuming that all the eigenvalues have multiplicity one, and we optimize iteratively over each unit vector q_i . We start by q_1 : augmenting the Lagrangian with the constraint $\|q_1\|^2 = 1$ using the Lagrange multiplier λ_1 and differentiating, we obtain:

$$\nabla_{q_1}(AUC(Q) + \lambda\|q_1\|^2) = \Sigma q_1 + \lambda_1 q_1 = 0.$$

This means that q_1 is a unit eigenvector of Σ . Denoting $\pi(1)$ the index of this eigenvector, we have $q_1^* = u_{\pi(1)}$ and the eigenvalue is $-\lambda_1 = \sigma_{\pi(1)}^2$. The criterion for this partially optimal value is now

$$AUC([q_1^*, q_2 \dots q_n]) = \sigma_{\pi(1)}^2 + \sum_{i=2}^k i q_i^T \Sigma q_i + (k+1) \sum_{i=k+1}^n q_i^T \Sigma q_i.$$

To take into account the orthogonality of the remaining vectors q_i ($i > 1$) with q_1^* in the optimization, we can project all the above quantities along $u_{\pi(1)}$. Optimizing now for q_2 under the constraint $\|q_2\|^2 = 1$, we find that q_2 is a unit eigenvector of $\Sigma - \sigma_{\pi(1)}^2 u_{\pi(1)} u_{\pi(1)}^T$ associated to a non-zero eigenvalue. Denoting $\pi(2)$ the index of this eigenvector (which is thus different from $\pi(1)$ because it has to be non-zero), we have $q_2^* = u_{\pi(2)}$ and the eigenvalue is $-\lambda_2 = 2\sigma_{\pi(2)}^2$.

Iterating the process, we conclude that $q_i^* = u_{\pi(i)}$ for some permutation π of the indices $1, \dots, n$. Moreover, the value of the criterion for that permutation is

$$AUC([q_1^*, q_2^* \dots q_n^*]) = \sum_{i=q}^k i \sigma_{\pi(i)}^2 + (k+1) \sum_{i=k+1}^n \sigma_{\pi(i)}^2.$$

In order to find the global minimum, we now have to compare the values of this criterion for all the possible permutations.

Assuming that $i < j$, we now show that the permutation of two indices $\pi(i)$ and $\pi(j)$ give a lower (or equal) criterion when $\pi(i) < \pi(j)$. Because eigenvalues are sorted in strictly decreasing order, we have $\sigma_{\pi(i)}^2 > \sigma_{\pi(j)}^2$. Thus, $(\alpha-1)\sigma_{\pi(i)}^2 > (\alpha-1)\sigma_{\pi(j)}^2$ for any $\alpha \geq 1$ and adding $\sigma_{\pi(i)}^2 + \sigma_{\pi(j)}^2$ on both sides, we get $\alpha\sigma_{\pi(i)}^2 + \sigma_{\pi(j)}^2 > \sigma_{\pi(i)}^2 + \alpha\sigma_{\pi(j)}^2$. For the value of α , we distinguish there cases:

- $i < j \leq k$: we take $\alpha = j/i > 1$. multiplying on both sides by the positive value i , we get: $i\sigma_{\pi(i)}^2 + j\sigma_{\pi(j)}^2 < i\sigma_{\pi(j)}^2 + j\sigma_{\pi(i)}^2$. The value of the criterion is thus strictly lower if $\pi(i) < \pi(j)$.
- $i \leq k < j$: we take $\alpha = (k+1)/i > 1$ and we get: $i\sigma_{\pi(i)}^2 + (k+1)\sigma_{\pi(j)}^2 < i\sigma_{\pi(j)}^2 + (k+1)\sigma_{\pi(i)}^2$. Once again, the value of the criterion is thus strictly lower if $\pi(i) < \pi(j)$.
- $k < i < j$: here permuting the indices does not change the criterion since $\sigma_{\pi(i)}^2$ and $\sigma_{\pi(j)}^2$ are both counted with the weight $(k+1)$.

In all cases, the criterion is strictly minimized by swapping indices in the permutation such that $\pi(i) < \pi(j)$ for $i < j$ and $i < k$. The global minimum is thus achieved for the identity permutation $\pi(i) = i$ for the indices $1 \leq i \leq k$. For the higher indices, any linear combination of the last $n - k$ eigenvectors of Σ gives the same value of the criterion. Taking into account the orthonormality constraints, such a linear combination writes $[q_{k+1}, \dots, q_n] = [u_{k+1}, \dots, u_n]R$ for some orthonormal $(n - k) \times (n - k)$ matrix R .

When some eigenvalues of Σ have a multiplicity larger than one, then the corresponding eigenvectors cannot be uniquely determined since they can be rotated within the eigenspace. With our assumptions, this can only occur within the last $n - k$ eigenvalues and this does not change anyway the value of the criterion. We have thus proved the following theorem.

Theorem 9 (Euclidean PCA as an optimization in the flag space).

Let $\hat{Y} = \{\hat{y}_i\}_{i=1}^N$ be a set of N data points in \mathbb{R}^n . We denote as usual the mean by $\bar{y} = \frac{1}{N} \sum_{i=1}^N \hat{y}_i$ and the empirical covariance matrix by $\Sigma = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - \bar{y})(\hat{y}_i - \bar{y})^T$. Its spectral decomposition is denoted $\Sigma = \sum_{j=1}^n \sigma_j^2 u_j u_j^T$ with the eigenvalues sorted in decreasing order. We assume that the first $k+1$ eigenvalues have multiplicity one, so that the order from σ_1 to σ_{k+1} is strict.

Then the partial flag of affine subspaces $Fl(x_0 \prec x_1 \dots \prec x_k)$ optimizing the AUC criterion:

$$AUC(Fl(x_0 \prec x_1 \dots \prec x_k)) = \sum_{i=0}^k \sigma_{out}^2(Fl_i(x_0 \prec x_1 \dots \prec x_k))$$

is totally ordered and can be parameterized by $x_0 = \bar{y}$, $x_i = x_0 + u_i$ for $1 \leq i \leq k$. The parametrization by points is not unique but the flag of subspaces which is generated is and is equal to the flag generated by the PCA modes up to mode k included.

References

- P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2008.
- Bijam Afsari. Riemannian l^p center of mass: existence, uniqueness, and convexity. *Proc. of the AMS*, 180(2):655–673, February 2010.
- R. Bhattacharya and V. Patrangenaru. Large sample theory of intrinsic and extrinsic sample means on manifolds, I. *Annals of Statistics*, 31(1):1–29, 2003.
- R. Bhattacharya and V. Patrangenaru. Large sample theory of intrinsic and extrinsic sample means on manifolds, II. *Annals of Statistics*, 33(3):1225–1259, 2005.

- Leo Brewin. Riemann normal coordinate expansions using cadabra. *Classical and Quantum Gravity*, 26(17):175017, 2009. URL <http://iopscience.iop.org/0264-9381/26/17/175017>.
- Peter Buser and Hermann Karcher. *Gromov's almost flat manifolds*. Société mathématique de France, 1981.
- Samuel R. Buss and Jay P. Fillmore. Spherical Averages and Applications to Spherical Splines and Interpolation. *ACM Trans. Graph.*, 20(2):95–126, April 2001. ISSN 0730-0301. doi: 10.1145/502122.502124. URL <http://doi.acm.org/10.1145/502122.502124>.
- Peter Crouch and F. Silva Leite. The dynamic interpolation problem: on Riemannian manifolds, Lie groups, and symmetric spaces. *Journal of Dynamical and control systems*, 1(2):177–202, 1995. URL <http://link.springer.com/article/10.1007/BF02254638>.
- James Damon and J. S. Marron. Backwards Principal Component Analysis and Principal Nested Relations. *Journal of Mathematical Imaging and Vision*, 50(1-2):107–114, October 2013. ISSN 0924-9907, 1573-7683. doi: 10.1007/s10851-013-0463-2. URL <http://link.springer.com/article/10.1007/s10851-013-0463-2>.
- Benjamin Eltzner, Sungkyu Jung, and Stephan Huckemann. Dimension Reduction on Polyspheres with Application to Skeletal Representations. In Frank Nielsen and Frédéric Barbaresco, editors, *Geometric Science of Information*, number 9389 in Lecture Notes in Computer Science, pages 22–29. Springer International Publishing, October 2015. ISBN 978-3-319-25039-7 978-3-319-25040-3. URL http://link.springer.com/chapter/10.1007/978-3-319-25040-3_3. DOI: 10.1007/978-3-319-25040-3_3.
- P.T. Fletcher, C. Lu, S.M. Pizer, and S. Joshi. Principal geodesic analysis for the study of nonlinear statistics of shape. *IEEE Transactions on Medical Imaging*, 23(8):995–1005, August 2004. ISSN 0278-0062. doi: 10.1109/TMI.2004.831793.
- M. Fréchet. Les éléments aléatoires de nature quelconque dans un espace distancié. *Annales de l'Institut Henri Poincaré*, 10:215–310, 1948.
- F Gay-Balmaz, DD Holm, DM Meier, TS Ratiu, and F-X Vialard. Invariant higher-order variational problems. *Communications in Mathematical Physics*, 309:413–458, 2012. doi: 10.1007/s00220-011-1313-y. URL <http://dx.doi.org/10.1007/s00220-011-1313-y>.
- Alexander N. Gorban and Andrei Y. Zinovyev. Principal graphs and manifolds. In *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods and Techniques*, chapter 2, pages 28–59. 2009. doi: 10.4018/978-1-60566-766-9.

- David Groisser. Newton's method, zeroes of vector fields, and the Riemannian center of mass. *Adv. in Applied Math*, 33:95–135, 2004.
- Jacob Hinkle, P. Thomas Fletcher, and Sarang Joshi. Intrinsic Polynomials for Regression on Riemannian Manifolds. *Journal of Mathematical Imaging and Vision*, 50(1-2):32–52, February 2014. ISSN 0924-9907, 1573-7683. doi: 10.1007/s10851-013-0489-5. URL <http://link.springer.com/article/10.1007/s10851-013-0489-5>.
- S. Huckemann, T. Hotz, and A. Munk. Intrinsic shape analysis: Geodesic principal component analysis for Riemannian manifolds modulo Lie group actions. *Statistica Sinica*, 20:1–100, 2010.
- Stephan Huckemann and Herbert Ziezold. Principal component analysis for Riemannian manifolds, with an application to triangular shape spaces. *Advances in Applied Probability*, 38(2):299–319, June 2006. ISSN 0001-8678, 1475-6064. doi: 10.1239/aap/1151337073. URL <http://projecteuclid.org/euclid.aap/1151337073>.
- S. Jung, I. L. Dryden, and J. S. Marron. Analysis of principal nested spheres. *Biometrika*, 99(3):551–568, September 2012. ISSN 0006-3444, 1464-3510. doi: 10.1093/biomet/ass022. URL <http://biomet.oxfordjournals.org/cgi/doi/10.1093/biomet/ass022>.
- Sungkyu Jung, Xiaoxiao Liu, J. S. Marron, and Stephen M. Pizer. Generalized PCA via the Backward Stepwise Approach in Image Analysis. In Jorge Angeles, Benoit Boulet, James J. Clark, Jzsef Kvecses, and Kaleem Siddiqi, editors, *Proc. of the Int. Symposium Brain, Body and Machine*, number 83 in Advances in Intelligent and Soft Computing, pages 111–123. Springer Berlin Heidelberg, 2010. ISBN 978-3-642-16258-9, 978-3-642-16259-6. URL http://link.springer.com/chapter/10.1007/978-3-642-16259-6_9.
- H. Karcher. Riemannian center of mass and mollifier smoothing. *Communications in Pure and Applied Mathematics*, 30:509–541, 1977.
- Hermann Karcher. Riemannian Center of Mass and so called karcher mean. *arXiv:1407.2087 [math]*, July 2014. URL <http://arxiv.org/abs/1407.2087>. arXiv: 1407.2087.
- W.S. Kendall. Probability, convexity, and harmonic maps with small image I: uniqueness and fine existence. *Proc. London Math. Soc.*, 61(2):371–406, 1990.
- Huiling Le. Estimation of Riemannian barycenters. *LMS J. Comput. Math.*, 7: 193–200, 2004.
- John M. Lee. *Riemannian Manifolds: An Introduction to Curvature*. Springer, 1997.

- Marco Lorenzi and Xavier Pennec. Geodesics, Parallel Transport & One-parameter Subgroups for Diffeomorphic Image Registration. *International Journal of Computer Vision*, 105(2):111–127, November 2013. doi: 10.1007/s11263-012-0598-4. URL <https://hal.inria.fr/hal-00813835>.
- Marco Lorenzi, Nicholas Ayache, and Xavier Pennec. Regional flux analysis for discovering and quantifying anatomical changes: An application to the brain morphometry in Alzheimer’s disease. *NeuroImage*, 115:224–234, July 2015. doi: 10.1016/j.neuroimage.2015.04.051.
- L. Machado, F. Silva Leite, and K. Krakowski. Higher-order smoothing splines versus least squares problems on Riemannian manifolds. *Journal of Dynamical and Control Systems*, 16(1):121–148, January 2010. ISSN 1079-2724, 1573-8698. doi: 10.1007/s10883-010-9080-1. URL <http://link.springer.com/10.1007/s10883-010-9080-1>.
- Peter J. Olver. Geometric Foundations of Numerical Algorithms and Symmetry. *Applicable Algebra in Engineering, Communication and Computing*, 11(5):417–436, April 2001. ISSN 0938-1279, 1432-0622. doi: 10.1007/s002000000053. URL <http://link.springer.com/article/10.1007/s002000000053>.
- Xavier Pennec. Intrinsic Statistics on Riemannian Manifolds: Basic Tools for Geometric Measurements. *Journal of Mathematical Imaging and Vision*, 25(1):127–154, 2006. doi: 10.1007/s10851-006-6228-4. URL <https://hal.inria.fr/inria-00614994>.
- Xavier Pennec. Barycentric Subspaces and Affine Spans in Manifolds. In *Geometric Science of Information GSI’2015*, Proceedings of Geometric Science of Information GSI’2015, Palaiseau, France, October 2015. URL <https://hal.inria.fr/hal-01164463>.
- Xavier Pennec and Vincent Arsigny. Exponential Barycenters of the Canonical Cartan Connection and Invariant Means on Lie Groups. In Frederic Barbaresco, Amit Mishra, and Frank Nielsen, editors, *Matrix Information Geometry*, pages 123–168. Springer, May 2012. doi: 10.1007/978-3-642-30232-9_7. URL <https://hal.inria.fr/hal-00699361>.
- Xavier Pennec, Pierre Fillard, and Nicholas Ayache. A Riemannian Framework for Tensor Computing. *International Journal of Computer Vision*, 66(1):41–66, 2006. doi: 10.1007/s11263-005-3222-z. URL <https://hal.inria.fr/inria-00614990>.
- O. Shoval, H. Sheftel, G. Shinar, Y. Hart, O. Ramote, A. Mayo, E. Dekel, K. Kavanagh, and U. Alon. Evolutionary Trade-Offs, Pareto Optimality, and the Geometry of Phenotype Space. *Science*, 336(6085):1157–1160, June 2012. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1217405. URL <http://www.sciencemag.org/content/336/6085/1157>.

- Nikhil Singh, François-Xavier Vialard, and Marc Niethammer. Splines for diffeomorphisms. *Medical Image Analysis*, 25(1):56–71, October 2015. ISSN 1361-8415. doi: 10.1016/j.media.2015.04.012. URL <http://www.medicalimageanalysisjournal.com/article/S1361841515000626/abstract>.
- S. Sommer, F. Lauze, and M. Nielsen. Optimization over geodesics for exact principal geodesic analysis. *Advances in Computational Mathematics*, 40(2):283–313, June 2013. ISSN 1019-7168, 1572-9044. doi: 10.1007/s10444-013-9308-1. URL <http://link.springer.com/article/10.1007/s10444-013-9308-1>.
- Stefan Sommer. Horizontal Dimensionality Reduction and Iterated Frame Bundle Development. In Frank Nielsen and Frédéric Barbaresco, editors, *Geometric Science of Information*, number 8085 in Lecture Notes in Computer Science, pages 76–83. Springer Berlin Heidelberg, 2013. ISBN 978-3-642-40019-3, 978-3-642-40020-9.
- Michael E. Tipping and Chris M. Bishop. Probabilistic Principal Component Analysis. *Journal of the Royal Statistical Society, Series B*, 61:611–622, 1999.
- Cédric Villani. Regularity of optimal transport and cut locus: From nonsmooth analysis to geometry to smooth analysis. *Discrete and Continuous Dynamical Systems*, 30(2):559–571, February 2011. ISSN 1078-0947. doi: 10.3934/dcds.2011.30.559. URL <http://www.aims sciences.org/journals/displayArticles.jsp?paperID=5986>.
- Grady S. Weyenberg. *Statistics in the Billera-Holmes-Vogtmann treespace*. PhD thesis, University of Kentucky, 2015. URL http://uknowledge.uky.edu/statistics_etds/12.
- Le Yang. *Medians of probability measures in Riemannian manifolds and applications to radar target detection*. PhD thesis, Poitiers University, December 2011.
- Miaomiao Zhang and P. Thomas Fletcher. Probabilistic principal geodesic analysis. In *Advances in Neural Information Processing Systems*, pages 1178–1186, 2013. URL <http://papers.nips.cc/paper/5133-probabilistic-principal-geodesic-analysis>.