Principal Sub-manifolds

Zhigang Yao

Department of Statistics and Applied Probability National University of Singapore



Collaborator: Tung Pham (U of Melbourne)

・ロット (日) ・ (日) ・ (日)

3

Sar

Data on manifolds¹ may arise in (at least) two ways:

(1) Manifold is actual physical space where data reside
 → Usually sphere; from geophysics to marine biology



¹contrast to manifolds in data analysis, or manifold±valued data → < = → = → <

Data on Manifolds

Data on manifolds may arise in (at least) two ways:

- (2) Multivariate data under non-linear constraints, thus being forced onto manifold
 - $\hookrightarrow\,$ e.g. cones for positive-def matrices or Stiefel manifolds for ordered bases



A B K A B K

Data on Manifolds

$$D = \left(egin{array}{cccc} D_{11} & D_{12} & D_{13} \ D_{21} & D_{22} & D_{23} \ D_{31} & D_{32} & D_{33} \end{array}
ight): v^{ op}Dv > 0, v \in \mathbb{R}^3$$

- ► let P₃ be a space of all symmetric positive definite 3 by 3 matrices
- $\mathcal{P}_3 \subset \mathbb{R}^6$
- \mathcal{P}_3 is convex but not linear in \mathbb{R}^6 : $\frac{3}{2}D_1 \frac{1}{2}D_2$ might not in \mathcal{P}_3
- Its actual space is a cone in \mathbb{R}^4

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ○ ○ ○

Data on Manifolds

A k-ad (k landmarks)

$$\mathbf{z} = \{(\mathit{x_j}, \mathit{y_j}), 1 \leq j \leq k\}$$

Kendall's shape space Σ_2^k

$$\mathbf{z} - \langle \mathbf{z} \rangle \in \mathbb{R}^{2k-2} = \{(x_j, y_j)_{1 \le j \le k} : \Sigma x_j = 0, \Sigma y_j = 0\}$$

$$preshape : \mathbf{w} = \frac{\mathbf{z} - \langle \mathbf{z} \rangle}{|\mathbf{z} - \langle \mathbf{z} \rangle|} \in S^{2k-3} \subset \mathbb{R}^{2k-2} \text{(unit sphere}^2\text{)}$$

$$shape : e^{i\theta} \mathbf{w} \in S^{2k-3}/S^1$$

Zhigang Yao

Principal Flows on Manifolds

Э

DQC

Unexpected challenges: statistics more rooted in linearity than may originally meet the eye.

Even the simplest statistical notion (an average) is a highly non-trivial object.

May admit several "canonical" definitions, some of which do not guarantee existence, others may not guarantee uniqueness

e.g. Fréchet mean μ for random variable X on \mathcal{M}

$$\mathbb{E} d^2(\mu,X) \leq \mathbb{E} d^2(x,X) \qquad orall x \in \mathcal{M}.$$

・回 ・ ・ ヨ ・ ・ ヨ ・

Two Toy Data Structures

"S-shape" and "C-shape":



문 문 문

DQC

exp and log map

Exponential map: tangent space \mapsto manifold

ie exp
$$_{ar{x}}(\stackrel{
ightarrow}{ar{x}y})=y$$

Log map: manifold \mapsto tangent space

ie log
$$_{ar{x}}(y)=\stackrel{
ightarrow}{ar{x}y}$$



Principal Nested Spheres

1. Principal nested spheres (spherical case) paralleling nested subspaces (Jung/Dryden/Marron, 2012)

$$\mathbb{U}_1 \subset \mathbb{U}_2 \subset \cdots \mathbb{U}_{d-1} \subset \mathbb{U}_d = \mathbb{S}^d$$
 $\mathbb{U}_k = {
m arg min}_{\mathbb{S}^k \subset \mathbb{U}_{k+1}} \sum_{i=1}^n d_m^2(x_i, \mathbb{S}^k)$



프 🖌 🔺 프 🕨

A view on Tangent PCA³

- 1. Idea: local tangent space PCA at every point yields direction of local maximal variation at that point
 - \blacktriangleright For any $x \in \mathcal{M}$, take tangent space T_x , 'lift' data, take

$$\Sigma_h(ar{x}) = rac{1}{\sum_i \kappa_h\left(x_i,ar{x}
ight)} \sum_i \Big(\log_{ar{x}}(x_i) \otimes \log_{ar{x}}(x_i)\Big) \kappa_h\left(x_i,ar{x}
ight).$$

- Consider (orientation) field {λ₁(x̄)e₁(x̄)}_{x∈M} of 1st eigenvectors times 1st eigenvalues of Σ_h(x̄).
- Aim to capture non-geodesic variation
- Can make this well-defined and smooth ("same direction")
- 2. Question: starting at a point (e.g. Fréchet mean) is there an integral curve of 'maximal variation' along this field?

³(will not discuss regularity conditions)

・ロト ・回ト ・ヨト ・ヨト

Given a staring point x̄, transform the orientation field
V(x̄) = {-λ₁(x̄)e₁(x̄), λ₁(x̄)e₁(x̄)}
from Σ_h(x) into a vector field W such that ∀x ∈ N(x̄)
Σ_h(x) W(x) = λ₁(x) W(x) (i.e., W(x) ∈ V(x))
At least locally, within an open neighborhood N, we can pick the eigenvectors {e₁(x) : x ∈ N} to be pointing in the same direction

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ○ ○ ○

A Principal Flow?

What does "curve of maximal variation" mean?

Would like a reasonably smooth curve $\gamma(x)$ whose derivative $\dot{\gamma}(x) \in T_x M$ at any point $x \in M$ is pprox tangent to $\lambda_1(x)e_1(x)$

...AND maximizes the work done by the field on a particle traveling along its path





・ロト ・回 ト ・ヨト ・ヨト

590



・ロト ・回 ト ・ヨト ・ヨト

590



・ロト ・回 ト ・ヨト ・ヨト

590



Zhigang Yao Principal Flows on Manifolds

・ロト ・回 ト ・ヨト ・ヨト

590

A Principal Flow (Panaretos/Pham/Yao JASA, 2014)

(mod technicalities)⁴ Curve with midpoint \bar{x} , maximizing

$$\int \left| \left\langle \dot{\gamma}(t), W(\gamma(t)) \right\rangle \right| dt$$

SubM $(A, v, \mathcal{M}) = \left\{ \gamma : [0, r] \to \mathcal{M}, \gamma \in C^2(\mathcal{M}), \gamma(s) \neq \gamma(s') \text{ for } s \neq s',
ight.$ $\gamma(0) = A, \dot{\gamma}(0) = v, \ell(\gamma[0, t]) = t \text{ for all } 0 \le t \le r \le 1
ight\}.$

- ► Answer: yes, reformulate to Euler-Lagrange equations
- ▶ ∃ unique solution under mild conditions on manifold+field
- Requires geodesics and second fundamental tensor
- Numerically Feasible for many "standard" manifolds
- ► Canonical: reduces to ordinary PCA in Euclidean spaces

Principal Curve (Hastie/Stuetzle (1989)) and Examples ⁵



Figure (3.2) Each point on a principal curve is the average of the points that project there.



⁵Thanks to Trevor Hastie for sharing the exaples) $< \square > < \square > < = > < = > = <math>$

Quick Illustration (varying scale parameter h)





→ 御 → → 注 → → 注 →

E

DQC

Simulation-Sphere (S-shape, C-shape, Diffusion)



Zhigang Yao Principal Flows on Manifolds

590

Some notations

For any point x in \mathcal{M} and \mathcal{N} ,

$$egin{aligned} \Sigma_{x,\mathcal{M}} &= rac{1}{n}\sum_{i=1}^n \log_{x,\mathcal{M}}(x_i)\otimes \log_{x,\mathcal{M}}(x_i) ext{ and } \ \Sigma_{x,\mathcal{N}} &= rac{1}{n}\sum_{i=1}^n \log_{x,\mathcal{N}}(x_i)\otimes \log_{x,\mathcal{N}}(x_i). \end{aligned}$$

- Let $\{e_1(x, \mathcal{M}), \ldots, e_k(x, \mathcal{M})\}$ be the first k eigenvectors of $\Sigma_{x,\mathcal{M}}$. k eigenvectors $\{e_1(x, \mathcal{N}), \ldots, e_k(x, \mathcal{N})\}$ for $\Sigma_{x,\mathcal{N}}$.
- ► Let H_k(x, M) be the hyperplane on M spanned by {e₁(x, M),..., e_k(x, M)}; H_k(x, N) be the hyperplane on N spanned by{e₁(x, N),..., e_k(x, N)}.

San

$$\mathbf{B}(x,\mathcal{N},\epsilon)=ig\{y\in\mathcal{N}:d_\mathcal{N}(x,y)\leq\epsilonig\}.$$

where $d_{\mathcal{N}}(x, y)$ is the distance of x and y on \mathcal{N} .

- For any positive integer number k < m, and any point x ∈ M, let SubM(x, ε, k, M) be the set of all k-dimensional sub-manifolds of B(x, M, ε).
- ► For a given k, the main idea of the principal sub-manifold is: at each point B of the sub-manifold N, the sub-manifold should be able to explain the manifold variation as much as possible.

・ロト ・回ト ・ヨト ・ヨト

An ideal principal sub-manifold (Yao/Pham, 2016)

 $(Ideal principal sub-manifold)^6 k$ -th dimensional principal sub-manifold

$$rg \sup_{\mathcal{N}\in ext{SubM}ig(A,\epsilon,k,\mathcal{M}ig)} \int_{B\in\mathcal{N}} \Big(\cos(lpha_B) imes \sum_{j=1}^k \lambda_j(B,\mathcal{M}) \Big) d\mu_\mathcal{N},$$

where $\mu_{\mathcal{N}}$ is the measure on \mathcal{N}

- ► To measure the degree of variation, we use the angle α_B between the two hyperplanes, H_k(B, M) and H_k(B, N).
- ► Theoretically, if α_B = 0 for every B, then H_k(B, M) = H_k(B, N). For general cases, one would hope α_B is as small as possible.

⁶(subject to modification)

How to get a concrete sub-manifold

- The idea is to apply the mapping to map N into a ball of radius ε in its tangent space, and then we can use the polar coordinates of that ball in the tangent space.
- ▶ Denote L(N, ε) = log_p(N) to be the image of the sub-manifold N at p under the logarithm map.

・回 ・ ・ ヨ ・ ・ ヨ ・

Principal sub-manifold (Yao/Pham, 2016)

(principal sub-manifoly, slightly changed) *k*-th dimensional principal sub-manifold

$$rg \sup_{\mathcal{N}\in \mathrm{SubM}ig(A,\epsilon,k,\mathcal{M}ig)} \int_{\mathbf{log}_A(B)\in \mathbf{log}_A(\mathcal{N})} \Big(\cos(lpha_B') imes \sum_{j=1}^\kappa \lambda_j(B,\mathcal{M})\Big) d\mu_k,$$

where μ_k is the measure on the ball of the k-dimensional space of radius ϵ .

- Let v_B be the tangent vector field of a geodesic curve on the sub-manifold N from A to B, denote α'_B to be the angle between v_B and hyperplane H_k(B, M)

nan

Principal flow and principal submanifold



・ロト ・回ト ・ヨト ・ヨト

Э

DQC

Principal sub-manifold in Euclidean Space

Theorem

Assume that $\mathcal{M} = \mathbb{R}^d$ then

$$\arg \sup_{\substack{\mathcal{N} \in \operatorname{SubM}(A,\epsilon,k,\mathcal{M})}} \int_{\log_A(B) \in \log_A(\mathcal{N})} \left(\cos(\alpha_B) \times \sum_{j=1}^{\kappa} \lambda_j(B,\mathcal{M}) \right) d\mu_k$$

= Hyperplane spanned by $\{e_1(A,\mathcal{M}), e_2(A,\mathcal{M}), \dots, e_k(A,\mathcal{M})\}.$

・ロト ・四ト ・ヨト ・ヨト

Э

Algorithm 1: two-dimensional principal sub-manifold

- 1. At a point A (mean), use the log map: $\log_A(x_i) = y_i$.
- Find the covariance matrix from y₁,..., y_n

$$\Sigma_A = (y_i - A)^T \times (y_i - A)$$

3. Let $e_1(A)$ and $e_2(A)$ be the first and second eigenvectors of Σ_A . Define

$$Z_{l} = \epsilon \times \left[\cos \left(2l\pi/180 \right) e_{1}(A) + \sin \left(2l\pi/180 \right) e_{2}(A) \right],$$

with l = 1, ..., 180.

- Use exponential map to map Z_l on the manifold so we get a set of new points exp_A(Z_l) = A_l.
- Assume that we stay at point A_{l,i}, we are going to find A_{l,i+1} (A₁₀ = A and A_{l,1} = A_l)
 5.1 Find Σ_{Al,i}.
 5.2 Find ε₁(A_{1i}) and ε₂(A_{1i}).
 - 5.3 Find $\log_{A_{l,i}}(A_{l,i-1}) = v_{l,i}$.
 - 5.4 Find

 $u_{l,i} = \left\langle v_{l,i}, e_1\left(A_{l,i}\right) \right\rangle \times \left. e_1\left(A_{l,i}\right) + \left\langle v_{l,i}, e_2\left(A_{l,i}\right) \right\rangle \times \left. e_2\left(A_{l,i}\right) \right\rangle$

where $\langle a, b \rangle = \sum_{i=1}^{n} a_i b_i$ with $a = (a_1, \ldots, a_n)$ and $b = (b_1, \ldots, b_n)$.

5.5 Find

$$r_{l,i} = \epsilon \times \frac{u_{l,i}}{\|u_{l,i}\|}.$$

5.6 Then

$$A_{l,i+1} = \exp_{A_{l,i}} \left(- r_{l,i} \right)$$

5.7 Stop at A_{l,i+1} when

$$\left\langle \log_{A_{l,i+1}}(A_{l,i}), \log_{A_{l,i+1}}(x_j) \right\rangle \geq 0.$$

for all $j = 1, \ldots, n$.

- For every l = 1,..., 180, connect A_{l,i} with A_{l,i+1} for all i we get a net of principal sub-manifold.
- Output: A_l for 1 ≤ l ≤ 180.

프 > 프

A view of the projected two-dimensional sub-manifold



DQC

Visualization of the projected sub-manifold for data on S^3



イロト イヨト イヨト イヨト

Э

DQC

Principal sub-manifolds for four sea wave sets of data with noise on S^3



Figure 6: Principal sub-manifolds (with superimposed principal directions) for four sea wave sets of data with noise on S^3 . (a) Principal sub-manifolds with no noise added. (b), (c) and (d) provide the same information for three different noise levels.

イロト イヨト イヨト イヨト

Principal variation of leaf growth-data



Figure 9: Leaf growth over a growing period of Clone 1 (a), Clone 2 (b), Clone 3 (c), and a reference tree (d). (a) Four landmarks on the leaf of clone 1 have been connected and represented by a polygon at each growing period (27 polygons totally); (b)-(d) provide the same information for Clone 2 (22 polygons), Clone 3 (24 polygons) and the reference tree (31 polygons).

90

Principal variation of leaf growth-result



Figure 12: Principal sub-manifolds of the leaf growth data. (a) First principal direction obtained from the combined leaves at breast height and the crown of the reference tree; (b) Second principal direction obtained from the combined leaves at breast height and the crown of the reference tree. (c)-(h) provide the same information for clone 1, 2 and 3.

・ロト ・回ト ・ヨト ・ヨト

Э

nan

- 1. A manifold extension of principal components (retain canonical interpretation+ allow for more flexible reduction of non-geodesic variation)
- 2. Study of notions of local covariance on manifolds?
 - Behaviour of theoretical version of Σ_h(x) as process over x or h or both?
- 3. Asymptotics for empirical flows/sub-manifolds?
- 4. More generally? Notions of covariance for manifolds?

(1日) (日) (日)

- ▶ Panaretos, V., Pham, T., and Yao, Z. (2014). Principal Flows. JASA.
- Yao, Z. and Pham, T. (2016). Principal Sub-manifolds (2016). Manuscript.
- Liu, H., Yao, Z., Leung, S. and Chan, T. F. (2016). A Level Set Based Variational Principal Flow Method for Nonparametric Dimension Reduction on Riemannian Manifolds. *Manuscript*.

San