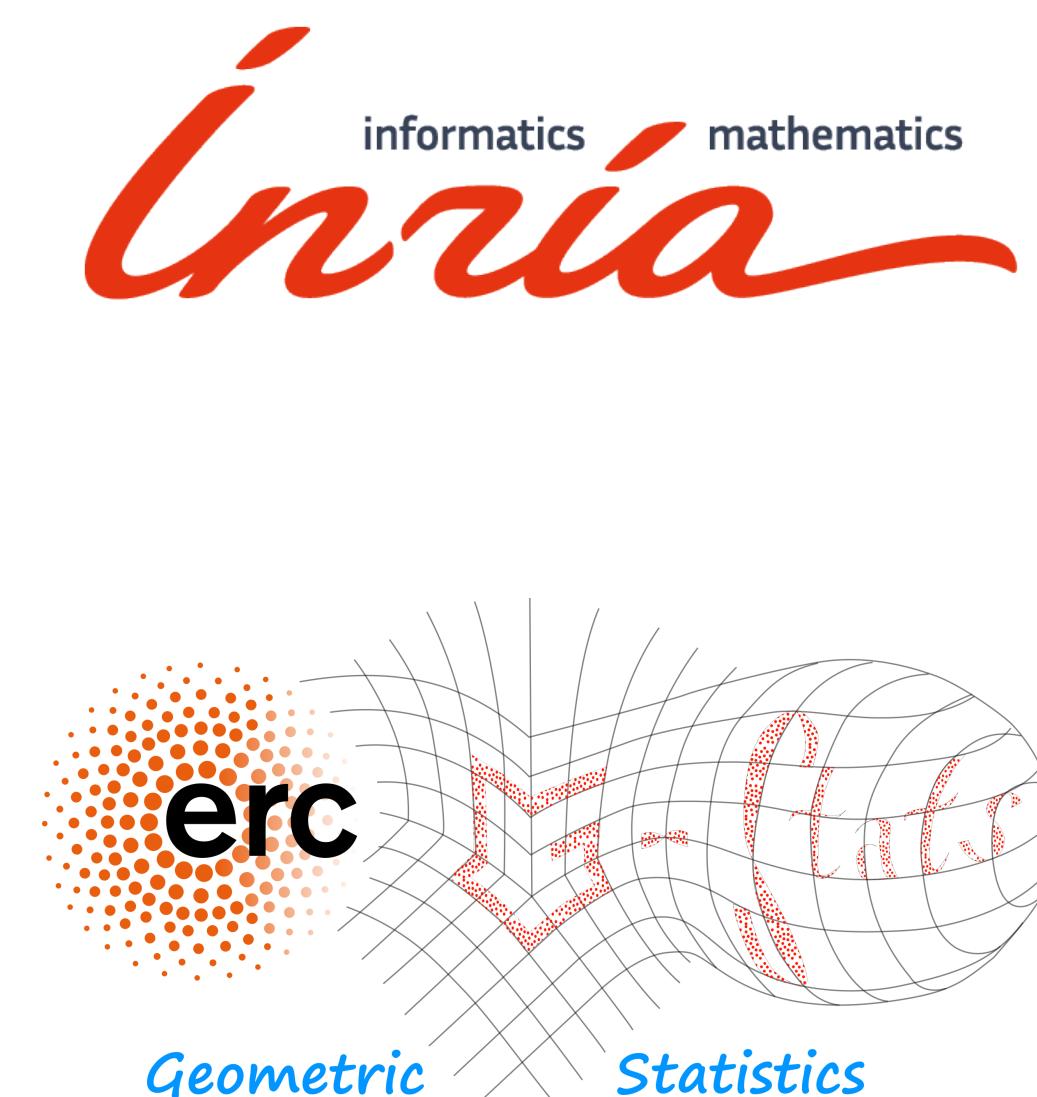


# Geomstats: A Python Package for Geometry in Machine Learning and Information Geometry

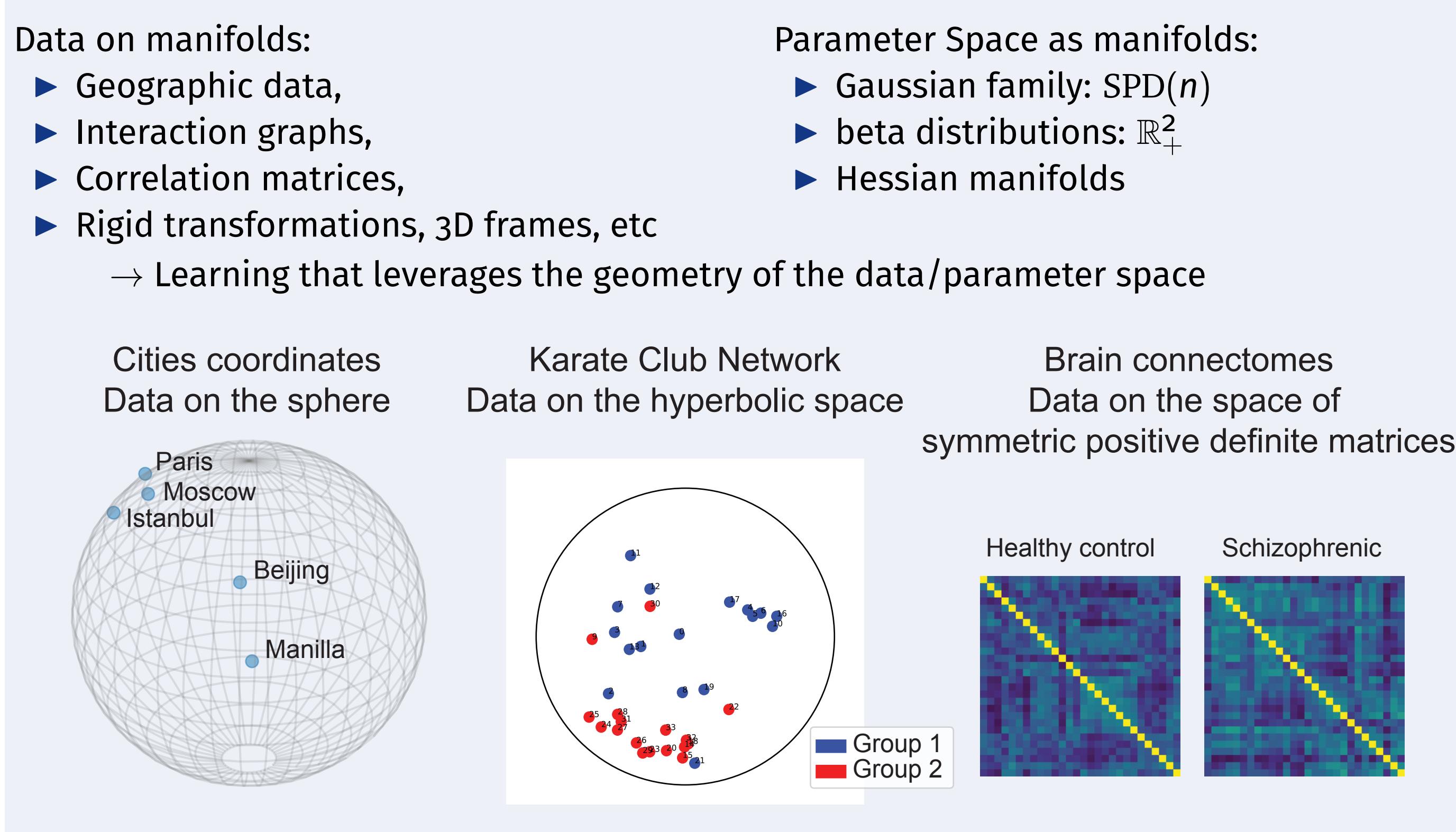
Nina Miolane, Nicolas Guigui<sup>1</sup>, Alice Le Brigant, Hadi Zaatiti, Christian Shewmake, Hatem Hajri, Johan Mathe, Benjamin Hou, Yann Thanwerdas, Stefan Heyder, Olivier Peltre, Niklas Koep, Yann Cabanes, Thomas Gerald, Paul Chauchat, Daniel Brooks, Bernhard Kainz, Claire Donnat, Susan Holmes, Xavier Pennec

<sup>1</sup>Université Côte d'Azur, Inria Sophia Antipolis, Epione Research Project, France.

Keywords: Riemannian Geometry, Python, Open-Source



## Motivation



## Objectives

- Teach "hands-on" Geometry, Learning and Information Geometry
- Democratize the use of Geometric Learning and Information Geometry in applications
- Support research in Geometric Learning and Information Geometry

## Package Description

Geomstats is an open-source Python package for computations and statistics on nonlinear manifolds [Miolane et al., 2020]. geomstats relies on three different back-ends, numpy, pytorch and tensorflow with a generic common API

```
export GEOMSTATS_BACKEND=numpy
import geomstats.backend as gs
```

The package implements tools for over 15 manifolds, each endowed with one or many Riemannian metrics.

- Spaces of constant curvature (Hypersphere, Hyperbolic space with different representations)
- Symmetric Positive Definite matrices
- Lie Groups (Rotations, Rigid-body transformations)
- Linear Subspaces (Grassmannian, Stiefel)
- Curve spaces (Discrete curves, Landmark spaces)
- Probability Distributions (Gaussians, Beta)
- Products of the above

Each module exposes methods to handle data that lie on the corresponding manifold, e.g. to sample random points, to check that a point indeed belongs to the manifold, to project vectors to tangent vectors to a point.

```
from geomstats.geometry.hypersphere import Hypersphere
sphere = Hypersphere(dim=2)
points = sphere.random_uniform(10)
print(sphere.belongs(points))
```

Riemannian Metrics expose methods to compute inner-products, geodesics, exponential, logarithm and parallel transport maps. Numerical methods are used when no closed-form solution exist

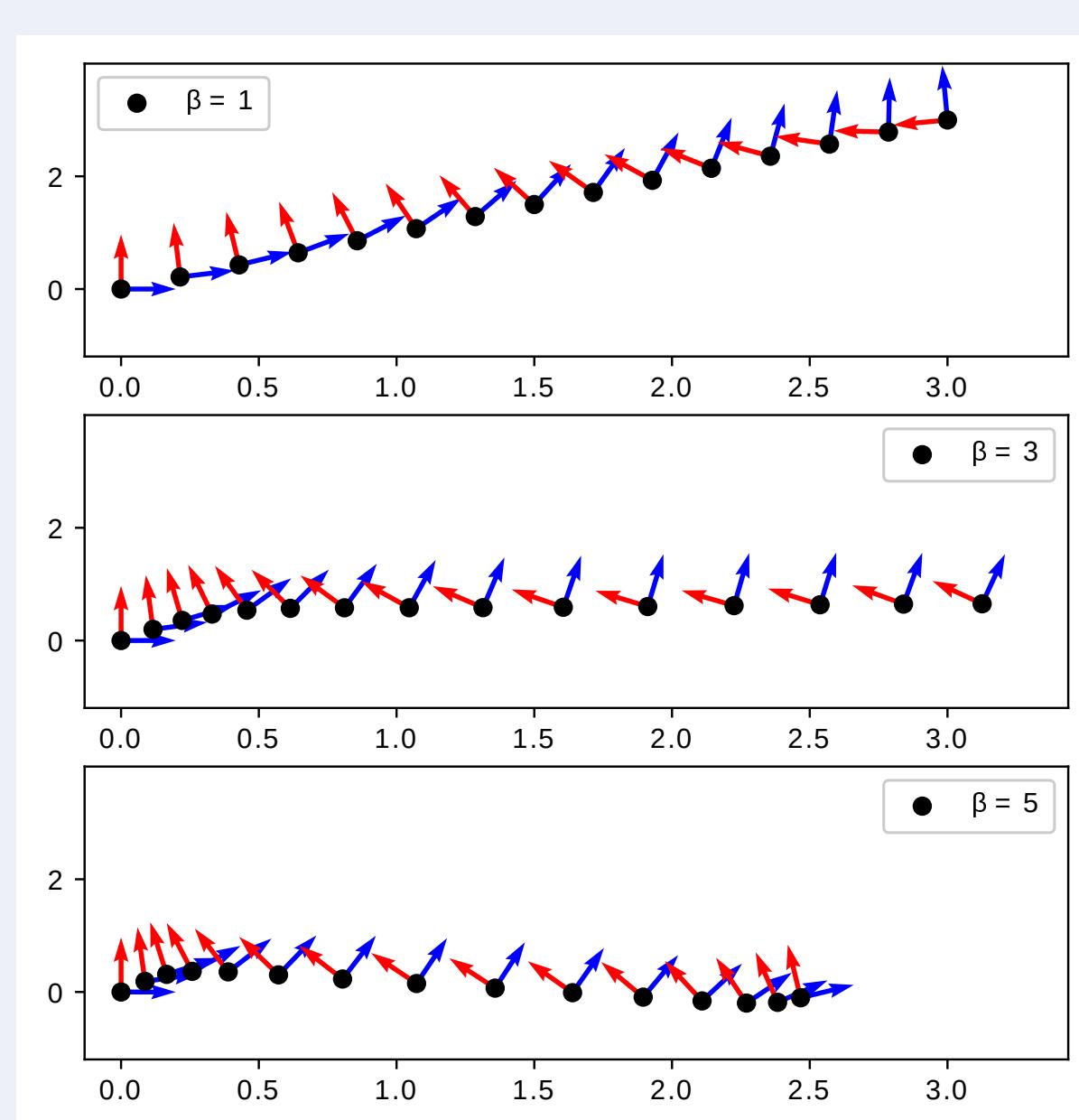


Figure 3: Geodesics in  $SE(2)$  with different left-invariant metrics.  $\beta$  is an anisotropy coefficient of the metric. Automatic differentiation and gradient descent allow to compute logarithms by geodesic shooting

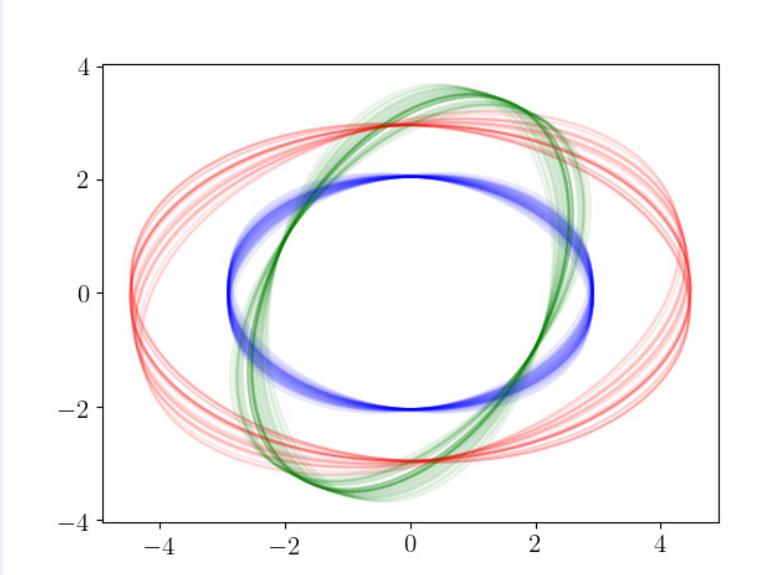


Figure 4: Geodesic balls on the manifold of Beta distributions with the Fisher metric

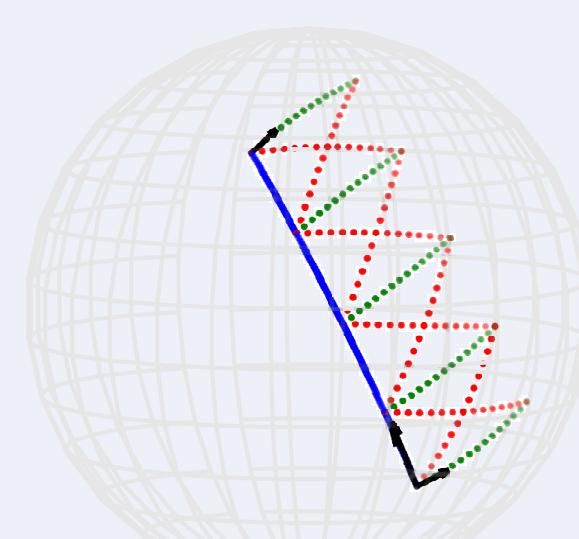
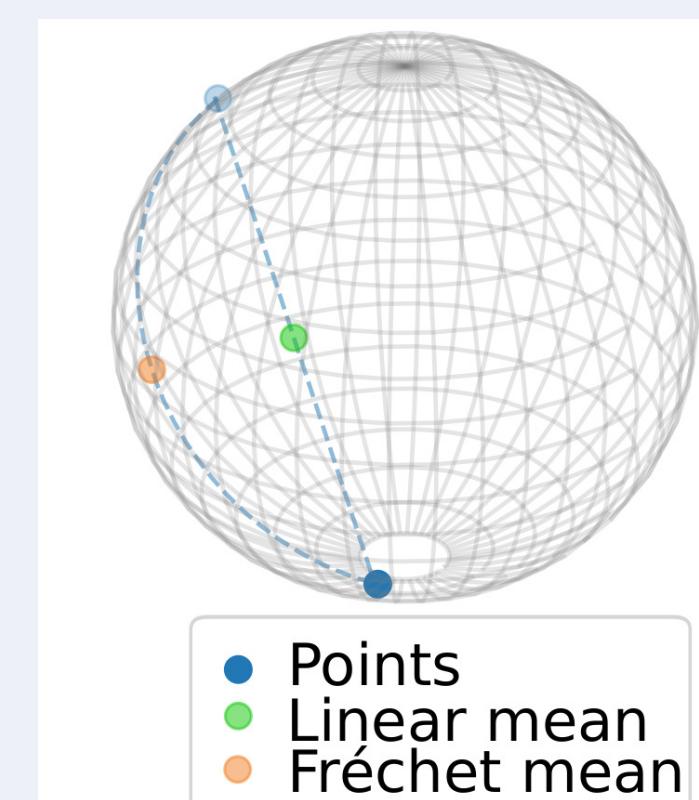


Figure 5: Schild's ladder: a numerical scheme to approach parallel transport [Guigui and Pennec, 2020]

## Fréchet Mean

- Generalized definition of the mean:  $\bar{x} = \operatorname{argmin}_{x \in M} \sum_{i=1}^n \operatorname{dist}_M(x, x_i)^2$
- Fréchet mean belongs to the manifold.

```
1 from geomstats.learning.frechet_mean import FrechetMean
2
3 estimator = FrechetMean(metric=sphere.metric)
4 estimator.fit(points)
5 frechet_mean = estimator.estimate_
```



## From Logistic Regression to tangent Logistic Regression

Use scikit-learn on the tangent space at the Fréchet mean

```
1 from geomstats.learning.preprocessing import ToTangentSpace
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.pipeline import make_pipeline
4
5 data, patient_ids, labels = data_utils.load_connectomes()
6 X_train, X_test, y_train, y_test = train_test_split(data, labels, random_state=0)
7
8 metric_aff = SPDMetricAffine(n=28)
9 pipeline = make_pipeline(ToTangentSpace(geometry=metric_aff), LogisticRegression(C=2))
10 pipeline.fit(X_train, y_train)
11 pipeline.predict(X_test)
> [0 1 1 0 0 1 0 0 1 1 1 0 0 0 1 0 1 1 1 1 0 0 0] ...
```

## From k-means to Riemannian k-means

Use Geomstats for unsupervised learning. Below an example of clustering on data from community network, embedded in the hyperbolic space [Gerald et al., 2020].

```
1 from geomstats.learning.kmeans import RiemannianKMeans
2
3 poincare_ball = PoincareBall(2)
4 karate_graph = data_utils.load_karate_graph()
5
6 hyperbolic_embedding = HyperbolicEmbedding()
7 embeddings = hyperbolic_embedding.embed(karate_graph)
8
9 kmeans = RiemannianKMeans(
10     riemannian_metric=poincare_ball.metric, n_clusters=2,
11     mean_method='frechet-poincare-ball')
12
13 centroids = kmeans.fit(X=embeddings, max_iter=100)
14 labels = kmeans.predict(X=embeddings)
```

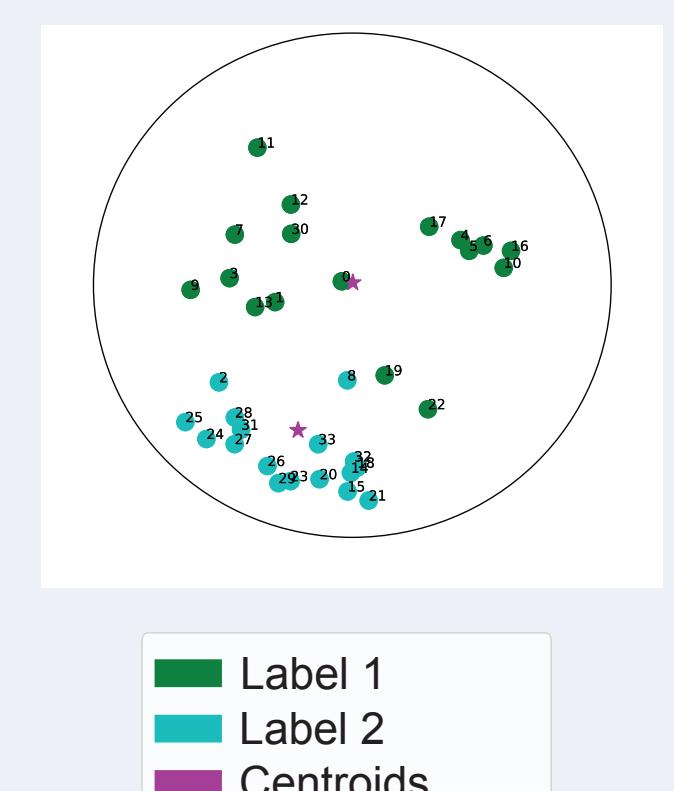


Figure 6: Riemannian k-means on the graph nodes after embedding in the Poincaré ball

## Beta distributions to classify histograms

We use maximum likelihood to fit beta distributions to histograms of cortical thickness maps [Brigant et al., 2020]. The manifold of beta distributions endowed with the Fisher-Rao metric has negative curvature. Geodesic distances can be computed numerically to apply k-nearest neighbor or k-means algorithms.

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from geomstats.geometry.beta_distributions import BetaDistributions
3
4 beta = BetaDistributions()
5 embeddings = beta.maximum_likelihood_fit(samples)
6
7 rnn = KNeighborsClassifier(n_neighbors=9, metric=beta.metric.dist)
8 rnn.fit(embeddings)
```

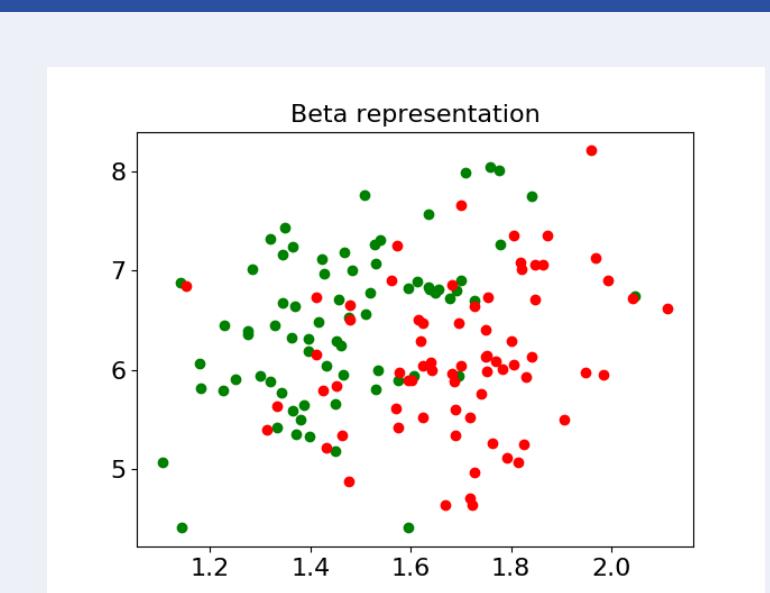


Figure 7: Embedding of the histograms in the 2D-manifold of beta distributions

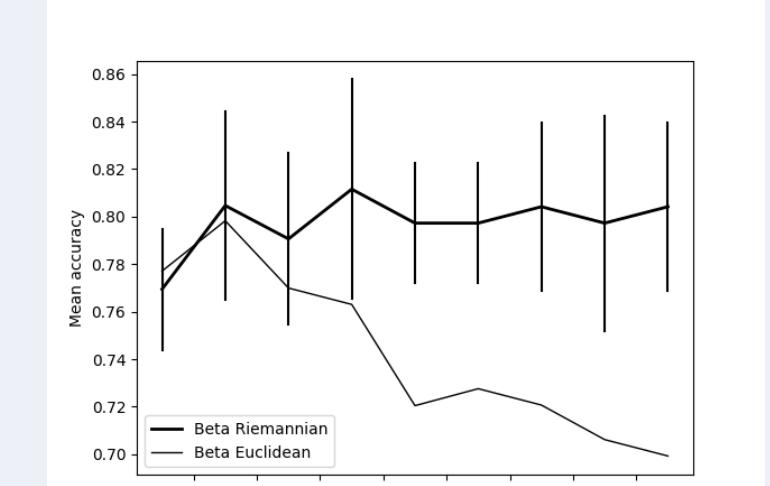


Figure 8: Mean accuracy of KNN classification on the cortical thickness data over 5-fold cross validation

## References

- Brigant, A. L., Guigui, N., Rebbah, S., and Puechmorel, S. (2020). Classifying histograms of medical data using information geometry of beta distributions.
- Gerald, T., Zaatiti, H., Hajri, H., Baskiotis, N., and Schwander, O. (2020). From Node Embedding To Community Embedding : A Hyperbolic Approach. arXiv: 1907.01662.
- Guigui, N. and Pennec, X. (2020). Numerical Accuracy of Ladder Schemes for Parallel Transport on Manifolds.
- Miolane, N., Brigant, A. L., Mathe, J., Hou, B., Guigui, N., Thanwerdas, Y., Heyder, S., Peltre, O., Koep, N., Zaatiti, H., Hajri, H., Cabanes, Y., Gerald, T., Chauchat, P., Shewmake, C., Kainz, B., Donnat, C., Holmes, S., and Pennec, X. (2020). Geomstats: A Python Package for Riemannian Geometry in Machine Learning.

