Introduction

Clustering algorithms

Resultats

Other distance

Discussion

# Improved Time-Series Clustering with UMAP dimension reduction method

C. Pealat, G. Bouleux, V. Cheutet

University of Lyon - DISP, INSA Lyon - University Hospital of Saint Etienne

2020



## Time series clustering

C. Pealat, G. Bouleux,V. Cheutet

### Introduction

Clustering algorithms

Resultats

Other distance

Discussion

Clustering : Find groups in unlabelled data



1-D Time series : specific data type

Objective : Determine the efficiency of UMAP as a pre-processing step for clustering algorithm.

## Introduction

Clustering algorithms

Resultats

Other distance

Discussion

# UMAP : Uniform Manifold Approximation and Projection

- Reduction of dimension algorithm, in general for visualization
- Methodology :

 $\hookrightarrow$  Creation of a fuzzy graph G : Edges between [0,1] with respect to the distance

 $\hookrightarrow$  Reduction of dimension of this graph with Laplacian Eigenmaps G'

 $\hookrightarrow$  Forced directed graph layout to minimized the entropy between G and G'

• What about clustering?

 $\hookrightarrow$  Controversy : It can create pseudo group

 $\hookrightarrow$  Test on real data

### Introduction

Clustering algorithms

Resultats

Other distance

Discussion

# Creation of a benchmark

- UCR Time Series Classification Archive
- 85 databases of time series with same length, with labels
- Benchmark of clustering results, using v-measure score to compare clustering results and true labels



Introduction

Clustering algorithms

Resultats

Other distance

Discussion

# Three Clustering algorithms

## • K-means

 $\hookrightarrow$  Centroid-based clustering

 $\hookrightarrow \underline{\text{Needs}}$ : Distance, definition of a mean, number of clusters<sup>1</sup>

## • Hierarchic

 $\hookrightarrow$  Connectivity-based clustering (Dendogram)

 $\hookrightarrow \underline{\text{Needs}} : Distance, number of clusters^1$ 

## • HDBSCAN

 $\hookrightarrow$  Density-based clustering

 $\hookrightarrow \underline{\mathrm{Needs}} : \mathit{Distance}$ 

<sup>1.</sup> Determined through silhouette score

# Results

#### <u>C. Pealat</u>, G. Bouleux,V. Cheutet

Introduction

Clustering algorithms

## Resultats

Other distance

Discussion

## <u>Results summarized :</u>

	Kmeans	Hier.	HDB.	U.+Kmeans	U.+Hier.	U.+HDB.
Mean	0.253	0.212	0.218	0.284	0.273	0.292
Best with U.				61%	66%	75%

 $\hookrightarrow$  With UMAP, improvement of the mean for all 3 algorithms

- $\hookrightarrow$  Improvement for at least 61% of the databases
- $\hookrightarrow$  Particularly efficient with HDBSCAN

And with an other distance?

Introduction

Clustering algorithms

Resultats

Other distance

Discussion

Results on the Stiefel manifold

Methodology new distance :

- ${\small \textcircled{0}}$  Embedding on  ${\bf R}^{{\bf n}\times {\bf p}}$  with delay coordinate embedding
- **2** Orthogonalization to obtain an element of  $V_{n,p} = \{A \in \mathbb{R}^{n \times p} : A^T A = Id\}$
- New similarity measure : geodesic on the Stiefel manifold (Principal Angle)

4 Karcher mean with respect to the geodesic, allows to apply K-means

	Kmeans	Hier.	HDB.	U.+Kmeans	U.+Hier.	U.+HDB.
Mean	0.203	0.173	0.182	0.292	0.213	0.301
Best with U.				50%	67%	83%

#### Introduction

Clustering algorithms

Resultats

Other distance

Discussion

# What to retain so

- Three clustering algorithms and two distances have been tested
- UMAP increased the v-measure score for all possible combinations
- In particular, UMAP coupled with HDBSCAN gave the best results
- UMAP and HDBSCAN can be still better fitted