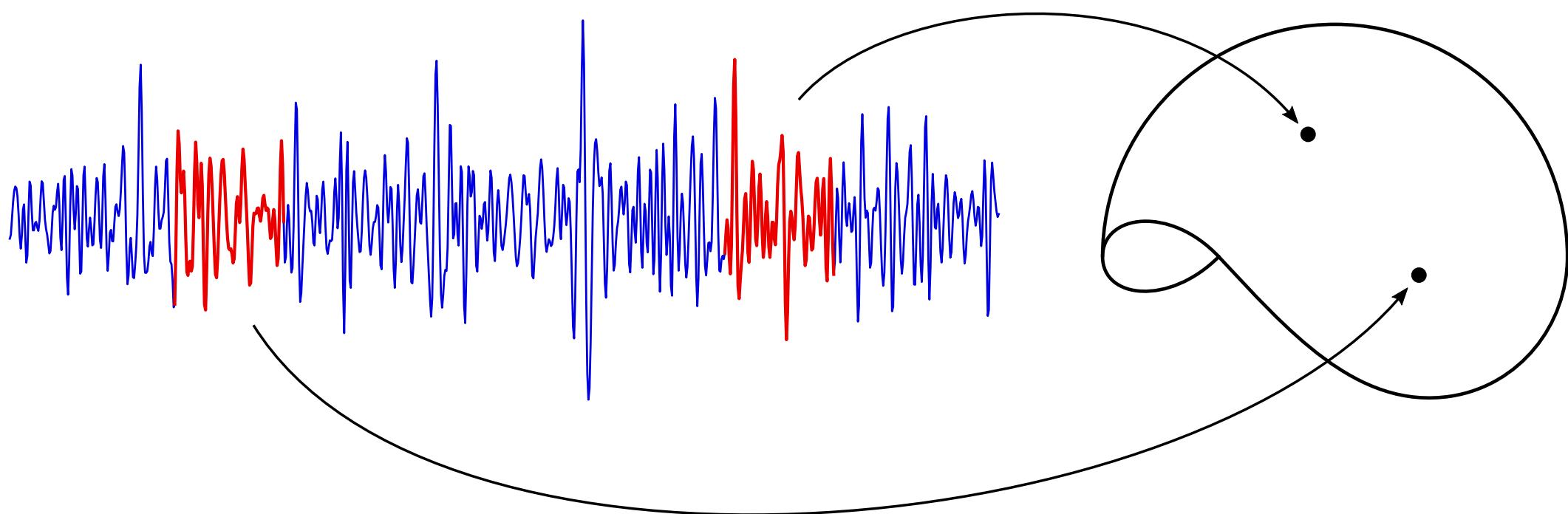


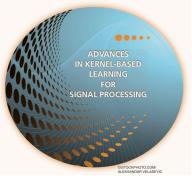
An introduction to diffusion maps and their use in time series analysis

Pedro Rodrigues – VIBS team meeting, 17th March 2017



Ronen Talmon, Israel Cohen, Sharon Gannot, and Ronald R. Coifman

Diffusion Maps for Signal Processing



A deeper look at manifold-learning techniques
based on kernels and graphs

Signal processing methods have significantly changed over the last several decades. Traditional methods were usually based on geometric statistical inference and linear filters. These frameworks have helped to solve a wide range of problems in signal processing, especially for implementation on digital signal processing (DSP) systems. Over the years, DSP systems have advanced rapidly, and their computational power has increased exponentially. Consequently, this development has enabled contemporary signal processing methods to handle much more complex data. Consequently, we have recently experienced a growing interest between signal processing and machine-learning approaches, e.g., dimensionality reduction, classification, and regression methods, whose computational burden is usually high.

In this article, we review manifold-learning techniques based on kernels and graphs. We present the main developments and trends and present ways to incorporate them into signal processing applications. We also discuss how to build compact representations of signals and intrinsic metrics and models, together with practical aspects and concrete signal

processing algorithms tackling challenging problems, e.g., transfer learning, inference, suppression and acoustic source localization. The prime focus is nonlinear signal processing using diffusion maps, which is a novel manifold-learning method. Our novel method and will encourage further research into the attractive and exciting new direction.

MOTIVATION AND BACKGROUND
In many applications in science and engineering, the observed data sets are controlled by underlying processes or drivers. As a result, these sets exhibit tightly redundant representations, and therefore, they can be compressed and represented more compactly and may enable a more efficient processing. For instance, molecular dynamics simulations of biologically significant macromolecules can be used to predict protein conformations for engineering and developing new drugs and treatments [1]; such simulations are often time-consuming and computationally expensive and often occur on time scales well beyond the computational reach of current solvers. Yet, by exploiting the (unknown) coherent structures in the data, the trajectories of the atoms can be more compactly represented by only a few, well-chosen rotation coordinates—for example, by a small subset of critical dihedral

Digital Object Identifier 10.1109/9.2012.2208050

Date of publication: 27 June 2012

1053-587X/13/070002-20\$31.00 © 2013 IEEE
IEEE SIGNAL PROCESSING MAGAZINE 79 JULY 2013

Ronen Talmon, Israel Cohen, Sharon Gannot, and Ronald Coifman (2013) **Diffusion maps for signal processing: A deeper look at manifold-learning techniques based on kernels and graphs**

Stéphane Lafon (Ph.D. thesis, 2004)
Diffusion maps and Geometric Harmonics

Ronald Coifman and Stéphane Lafon (2006)
Diffusion maps

Outline

Introduction

Diffusion maps

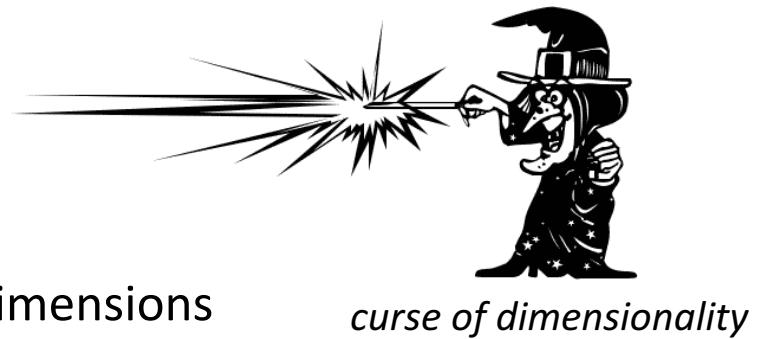
Time series analysis

Concluding remarks

Introduction

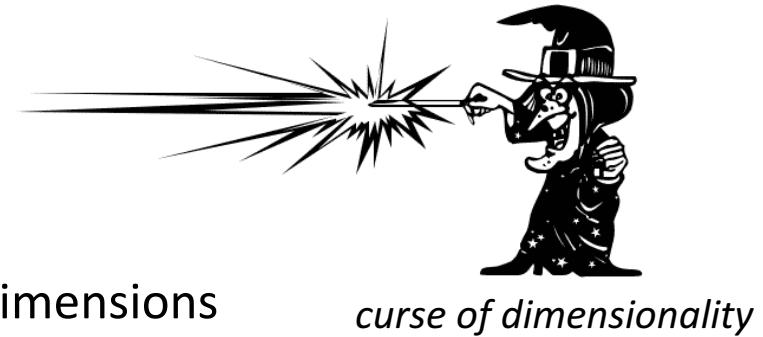
High-dimensionality is an obstacle to efficient data processing for a number of reasons

- Norms in \mathbf{R}^n are not **numerically equivalent** when n is large
- Density estimation is **difficult** in big dimensions (needs too much data)
- Fast algorithms in low dimensions may become **prohibitively slow** in high dimensions



High-dimensionality is an obstacle to efficient data processing for a number of reasons

- Norms in \mathbf{R}^n are not **numerically equivalent** when n is large
- Density estimation is **difficult** in big dimensions (needs too much data)
- Fast algorithms in low dimensions may become **prohibitively slow** in high dimensions

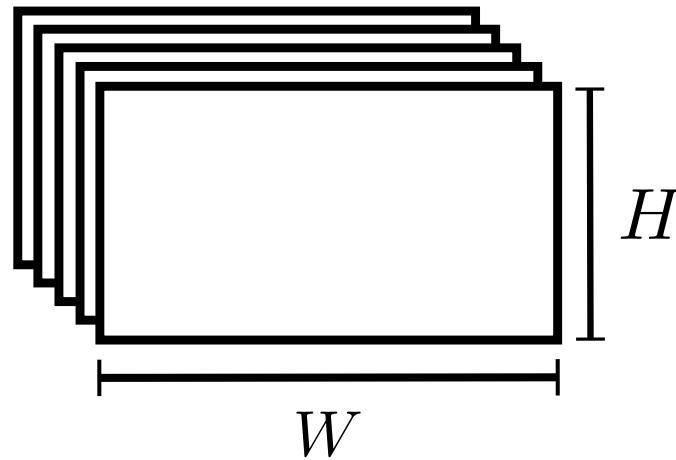


Data dimensions are often correlated via some **functional dependence**

Fundamental assumption: data lies on (or close to) a low-dimensional manifold

Dimension reduction techniques for visualization, feature selection, compression, etc.





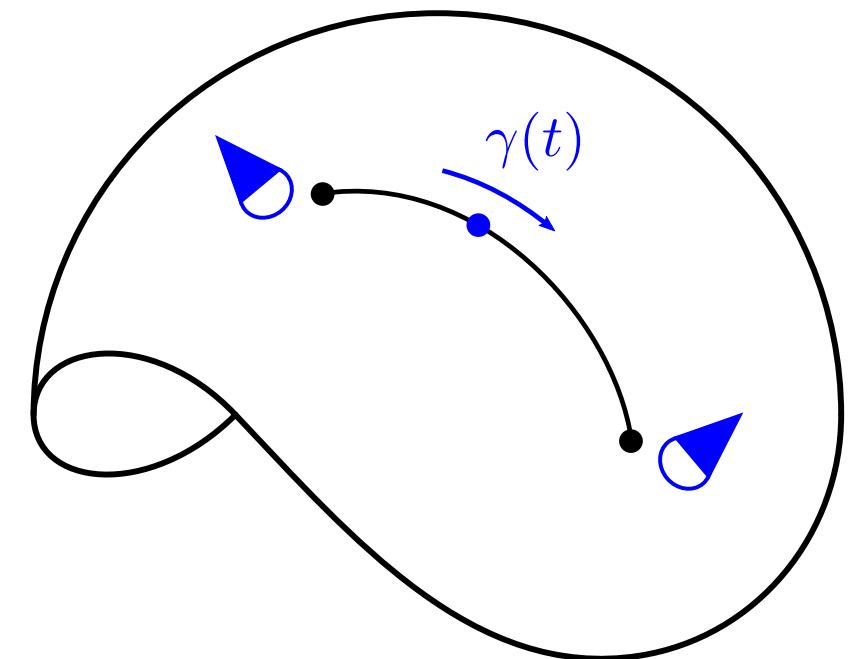
A video is composed of several **frames** $\left\{ X_i \right\}_{i=1}^N$

Each frame is an image and can be associated to a **vector**

Frames are points in a **very** high-dimensional space,

but there are not many **degrees of freedom** in the video

$$X_i \in \mathbb{R}^{W \times H} \Rightarrow \Phi(X_i) \in \mathbb{R}^d$$



$$X_i \in \mathbb{R}^n \Rightarrow \Phi(X_i) \in \mathbb{R}^d$$

PCA - Principal component analysis
MDS - Multi-dimensional scaling
LLE - Local-linear embedding

- Classical techniques like PCA and MDS for **linear** dimension reduction. Kernel PCA for **nonlinear**.
- Two papers in Science magazine (2000) presenting **manifold learning** techniques: Isomap and LLE
- Previous methods are all **global** (except LLE)
- Laplacian eigenmaps (2003) and Diffusion maps (2004) search for an embedding that respects **local** geometry

"Think globally, act locally"

Diffusion maps

The algorithm consists of **three** steps

- 1) Build a **graph** with samples as vertices and **weights** defined using a kernel
- 2) Define a **Markov process** on the graph (normalize **W** matrix)
- 3) Obtain a **nonlinear mapping** of the samples into a new embedded space

“A kernel specifies the **local geometry** of the data and captures some geometric features of interest. The Markov chain describes how this local information **propagates** along the graph”



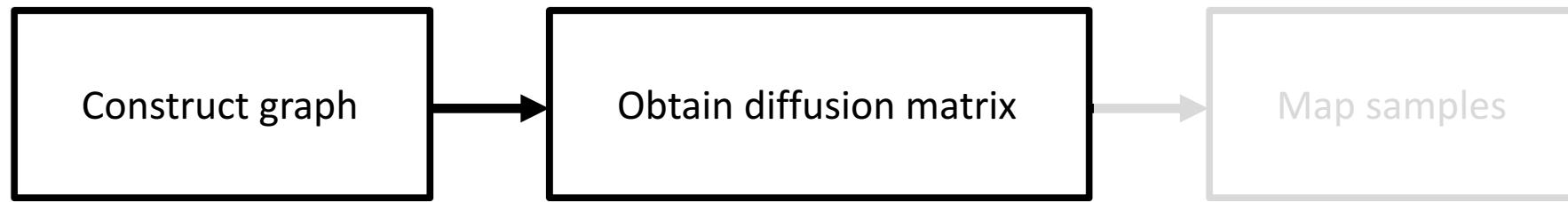
Define a graph Ω whose weights represent **pairwise affinity** between samples

$$\mathbf{W}(i, j) = \exp\left(-\frac{d(X_i, X_j)^2}{2\epsilon^2}\right)$$

ϵ scale parameter imposes
a notion of **locality**

$$d(X_i, X_j)$$

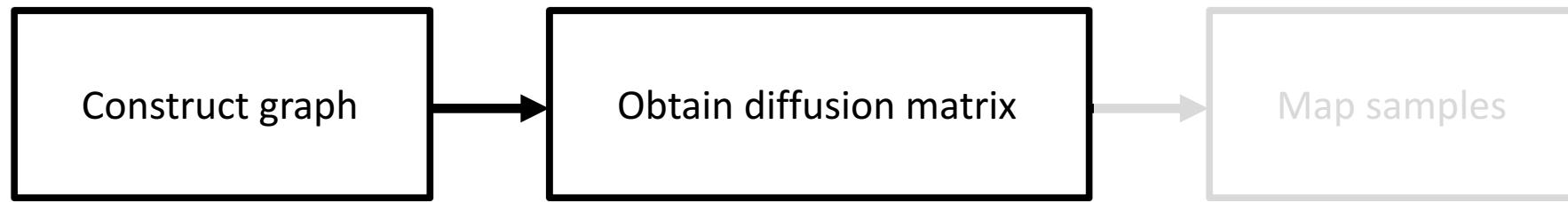
measure of **similarity**
between samples
(application-specific)



Normalize the rows of matrix \mathbf{W} so to have a probability transition (diffusion) matrix

$$\mathbf{P} = \mathbf{D}^{-1}\mathbf{W} \quad \mathbf{D}(i, i) = \sum_j \mathbf{W}(i, j)$$

$$\mathbf{P}(i, j) \geq 0 \quad \sum_j \mathbf{P}(i, j) = 1$$



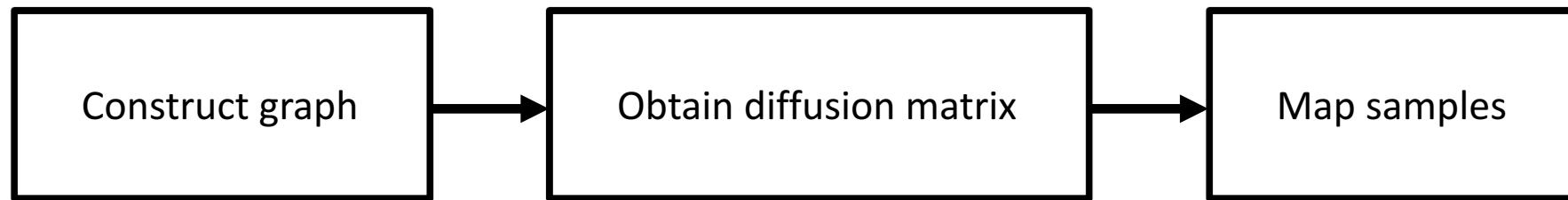
Normalize the rows of matrix \mathbf{W} so to have a probability transition (diffusion) matrix

$$\mathbf{P} = \mathbf{D}^{-1}\mathbf{W} \quad D(i, i) = \sum_j W(i, j)$$

probability of going from i to j in
a **random walk** over the graph
(diffusion process)

$$P(i, j) = \sum_{k \geq 0} \lambda_k \varphi_k(i) \psi_k(j)$$

$\{\psi_k\}$ **left** eigenvectors
 $\{\varphi_k\}$ **right** eigenvectors



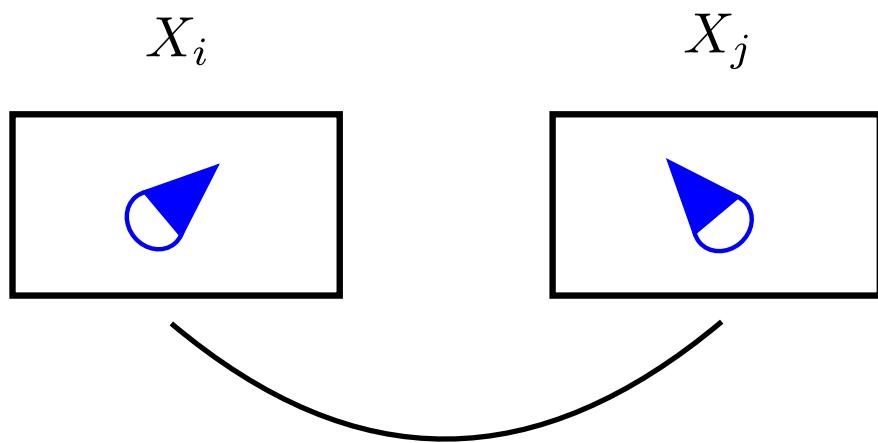
Check the **decay** of the eigenvalues and choose a dimension d

Embed the data points as **coordinates** of the right eigenvectors of \mathbf{P}

$$X_i \in \mathbb{R}^n \Rightarrow \Phi(X_i) = \begin{bmatrix} \lambda_1 \varphi_1(i) \\ \vdots \\ \lambda_d \varphi_d(i) \end{bmatrix} \in \mathbb{R}^d \quad \begin{aligned} \mathbf{P} \varphi_k &= \lambda_k \varphi_k \\ \varphi_k &\in \mathbb{R}^N \end{aligned}$$

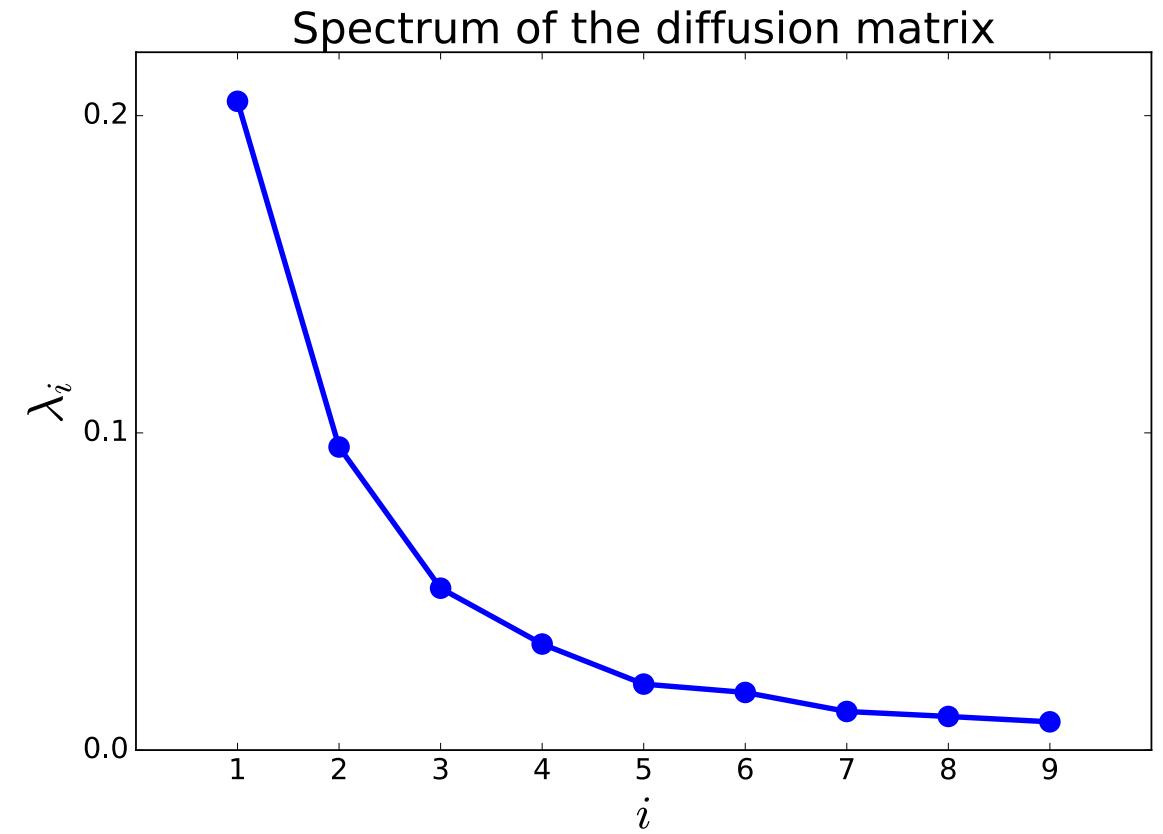
N is the number of samples

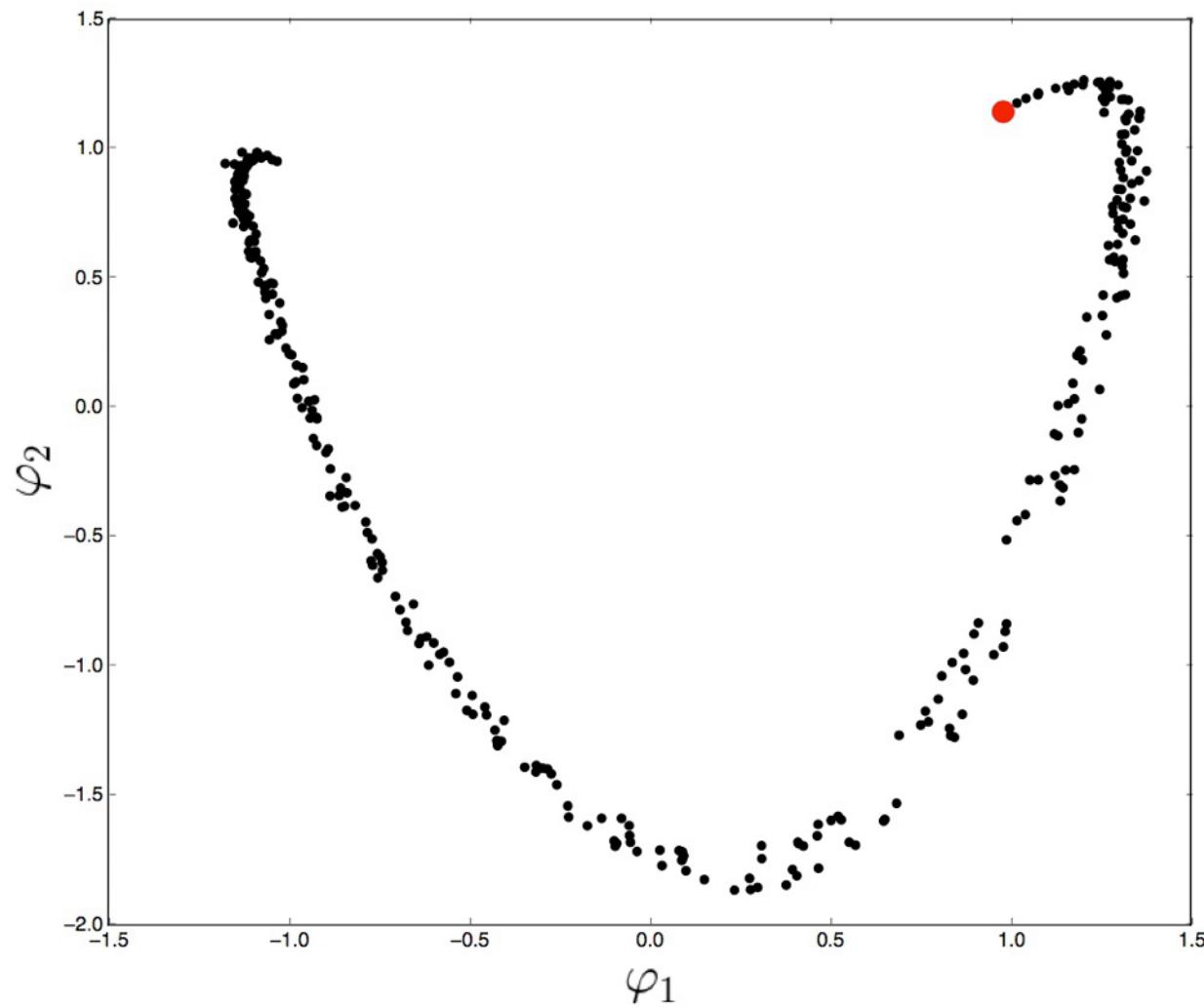
Example with **video** frames



$$d^2(X_i, X_j) = \sum_k \sum_{\ell} \left(X_i(k, \ell) - X_j(k, \ell) \right)^2$$

(other distances could be used)



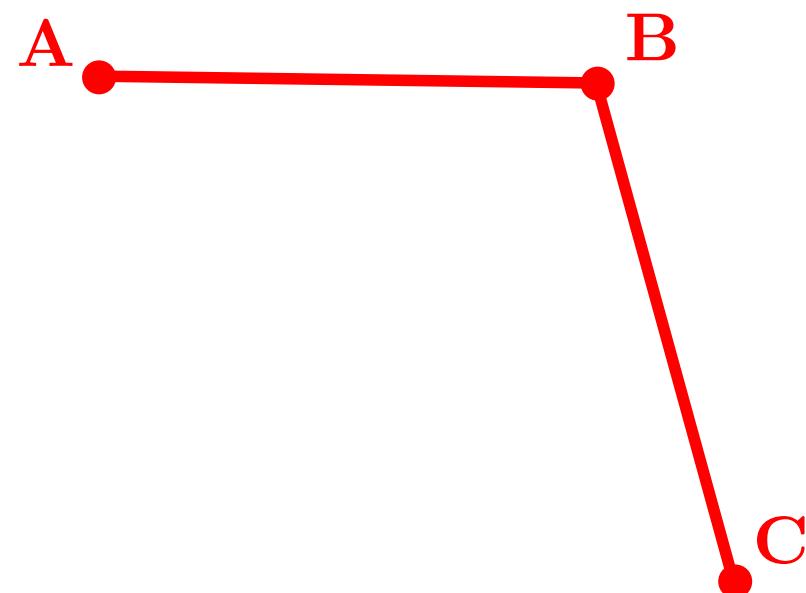


The embedding gives good results, but **why?**



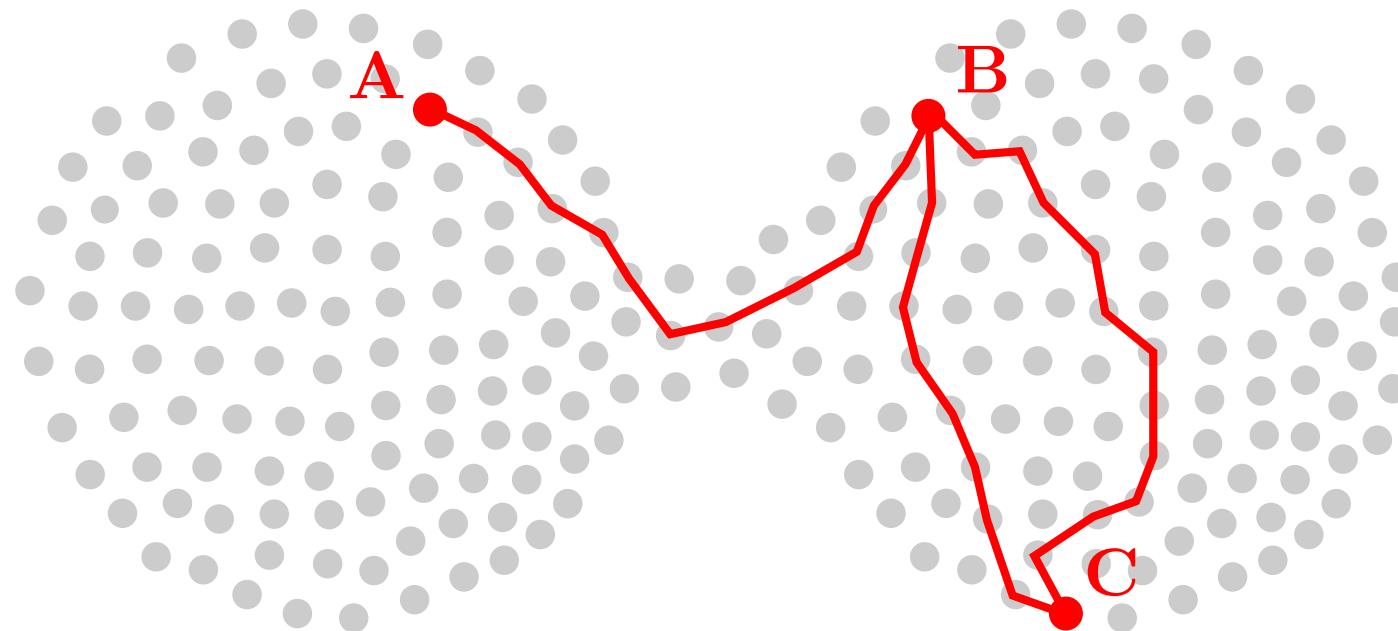
The embedding gives good results, but **why?**

It uses the notion of **diffusion distance**, a measure of connectivity between points in a dataset



The embedding gives good results, but **why**?

It uses the notion of **diffusion distance**, a measure of connectivity between points in a dataset



The embedding gives good results, but **why**?

It uses the notion of **diffusion distance**, a measure of connectivity between points in a dataset

$$D^2(X_i, X_j) = \sum_{y \in \Omega} \frac{1}{\pi(y)} \left(\mathbf{P}(i, y) - \mathbf{P}(j, y) \right)^2$$

which is the weighted L_2 distance between conditional distributions $\mathbf{P}(i, \cdot)$ and $\mathbf{P}(j, \cdot)$

$\pi(y)$ is the **stationary** distribution of the Markov chain defined on the graph Ω (first left eigenvector of \mathbf{P})

The embedding gives good results, but **why**?

It uses the notion of **diffusion distance**, a measure of connectivity between points in a dataset

$$D^2(X_i, X_j) = \sum_{y \in \Omega} \frac{1}{\pi(y)} \left(\mathbf{P}(i, y) - \mathbf{P}(j, y) \right)^2$$

- 1) Points are closer if they are highly **connected** in the graph (clusters)
- 2) Distances between two points depends on **all possible paths** between them (robust)
- 3) Distances between two points take into account **all evidences** relating them (good for inference)

Using the decomposition of the diffusion matrix

probability of going from i to j in
a **random walk** over the graph
(diffusion process)

$$\mathbf{P}(i, j) = \sum_{k \geq 0} \lambda_k \boldsymbol{\varphi}_k(i) \boldsymbol{\psi}_k(j)$$

$\{\boldsymbol{\psi}_k\}$ **left eigenvectors**
 $\{\boldsymbol{\varphi}_k\}$ **right eigenvectors**

we can show that

$$D^2(X_i, X_j) = \sum_{k \geq 0} \lambda_k^2 (\boldsymbol{\varphi}_k(i) - \boldsymbol{\varphi}_k(j))^2$$

and so

$$D^2(X_i, X_j) \simeq \sum_{k=1}^d \lambda_k^2 (\boldsymbol{\varphi}_k(i) - \boldsymbol{\varphi}_k(j))^2$$

$$\boldsymbol{\varphi}_0 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$D^2(X_i, X_j) \simeq \sum_{k=1}^d \lambda_k^2 \left(\varphi_k(i) - \varphi_k(j) \right)^2$$

this approximation can be interpreted as the **Euclidean** distance between two vectors

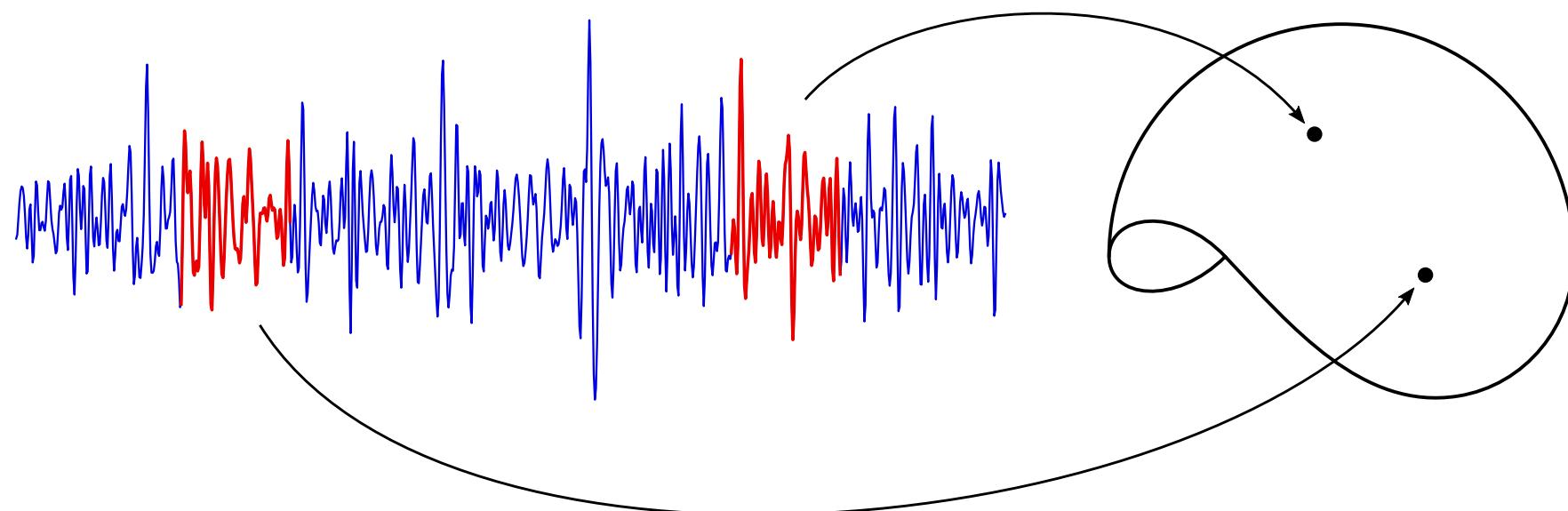
$$D^2(X_i, X_j) \simeq \|\Phi(X_i) - \Phi(X_j)\|^2$$

$$\Phi(X_i) = \begin{bmatrix} \lambda_1 \varphi_1(i) \\ \vdots \\ \lambda_d \varphi_d(i) \end{bmatrix} \quad \Phi(X_j) = \begin{bmatrix} \lambda_1 \varphi_1(j) \\ \vdots \\ \lambda_d \varphi_d(j) \end{bmatrix}$$

Time series analysis

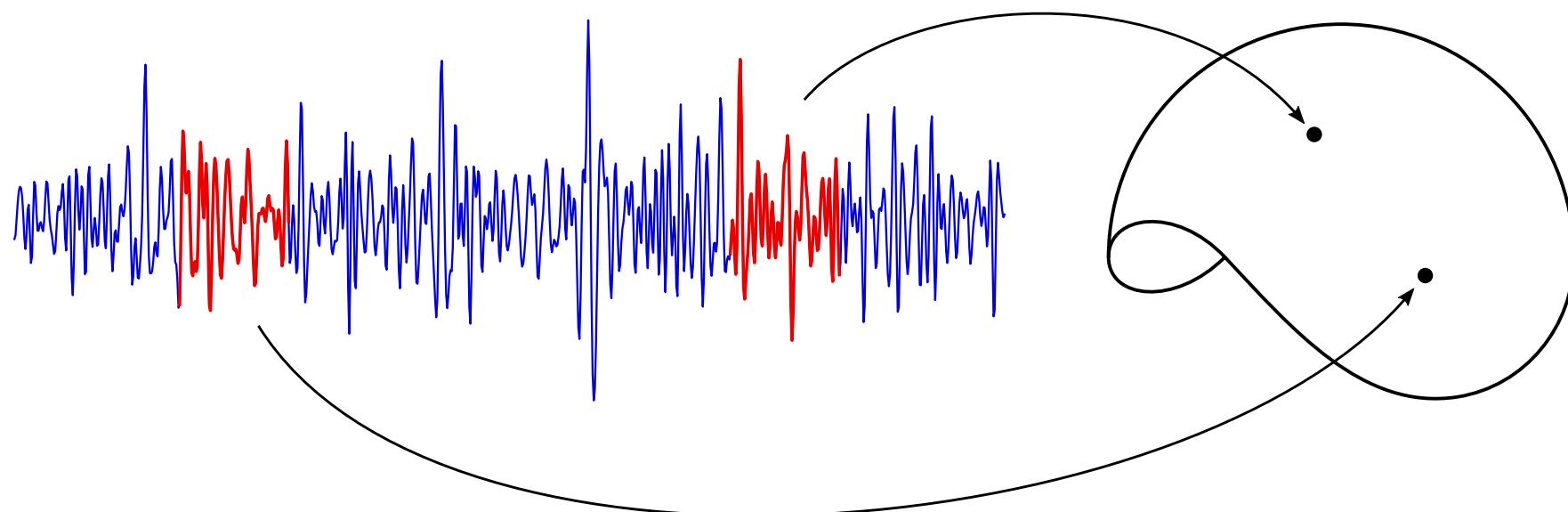
Two bottlenecks for time series analysis:

- Measurements are **realizations** of a stochastic process
- **Sequential** implementation



Two bottlenecks for time series analysis:

- Measurements are **realizations** of a stochastic process
- Sequential implementation



Describe signals in **state-space** and define an intrinsic affinity metric

$$\boldsymbol{x}(t) = f(\boldsymbol{\theta}(t))$$

The diagram illustrates the mapping from state to measurement. It features a central equation $\boldsymbol{x}(t) = f(\boldsymbol{\theta}(t))$. To the left of the equation, a horizontal line with a vertical tick at its left end points to the word "measurement". To the right, another horizontal line with a vertical tick at its right end points to the word "state". Below these, a third horizontal line with a vertical tick at its right end points to the text "map from state-space to measurement space (possibly nonlinear and/or probabilistic)".

measurement

state

map from state-space to measurement space
(possibly nonlinear and/or probabilistic)

Use an **observer operator** to describe the signals and analyze the geometry of the latent space

$$\boldsymbol{z}(t) = \mathcal{L}\{\boldsymbol{x}(t)\} \quad (\text{e.g. Fourier power spectral density})$$

Obtain the **statistics** of the observers in short windows (assuming local **stationarity**)

$$\hat{\boldsymbol{\mu}}(t) = \frac{1}{L_0} \sum_{\tau \in \mathcal{T}_t} \mathbf{z}(\tau) \quad (L_0 = |\mathcal{T}_t|)$$

$$\hat{\mathbf{C}}(t) = \frac{1}{L_0} \sum_{\tau \in \mathcal{T}_t} (\mathbf{z}(\tau) - \hat{\boldsymbol{\mu}}(t)) (\mathbf{z}(\tau) - \hat{\boldsymbol{\mu}}(t))^T$$

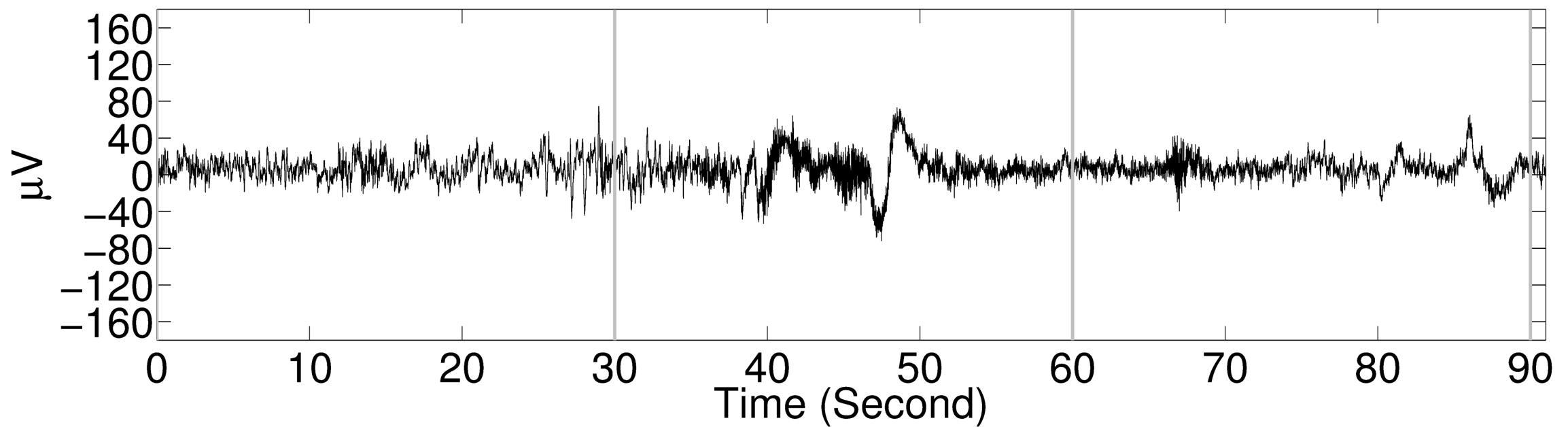
Define a **Mahalanobis** distance between samples

$$d^2(\mathbf{x}(t), \mathbf{x}(s)) = \frac{1}{2} (\mathbf{z}(t) - \mathbf{z}(s))^T (\mathbf{C}^\dagger(t) + \mathbf{C}^\dagger(s)) (\mathbf{z}(t) - \mathbf{z}(s))$$

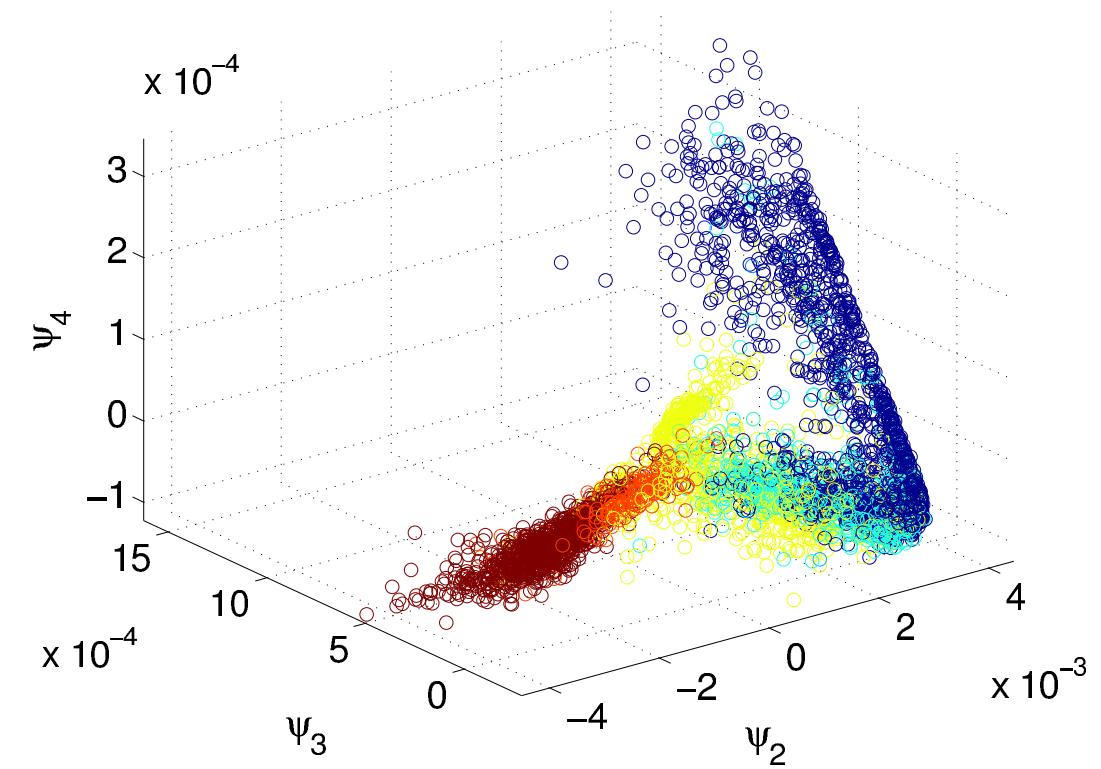
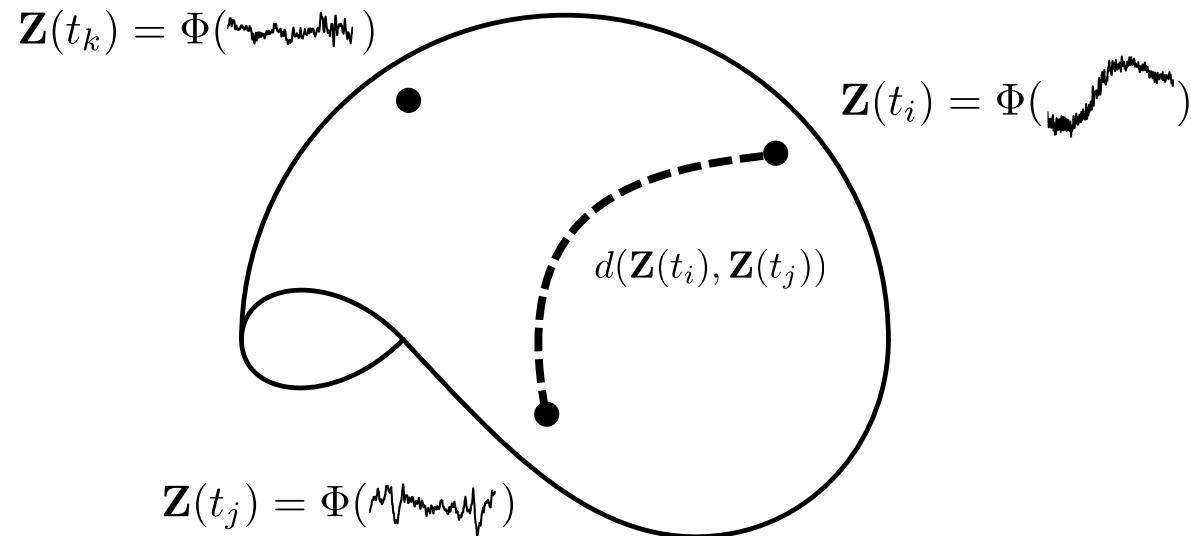
Hau-tieng Wu, Ronen Talmon, and Yu-Lun Lo

Assess sleep stage by modern signal processing techniques

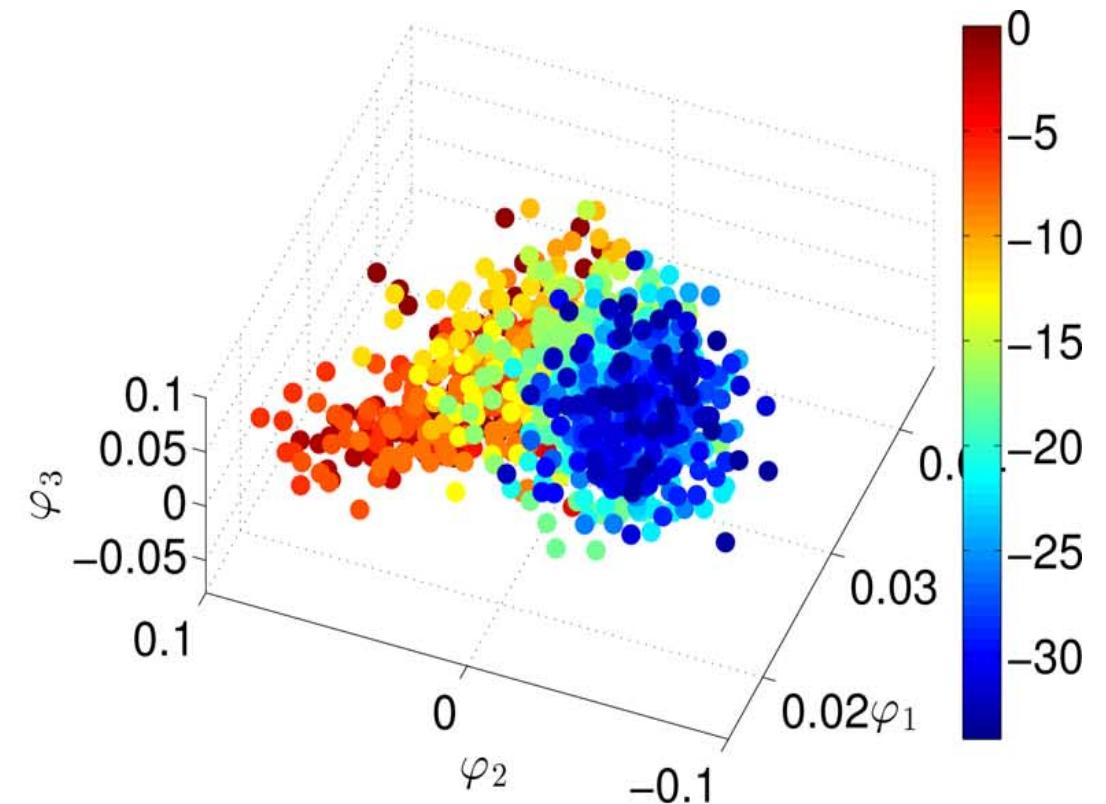
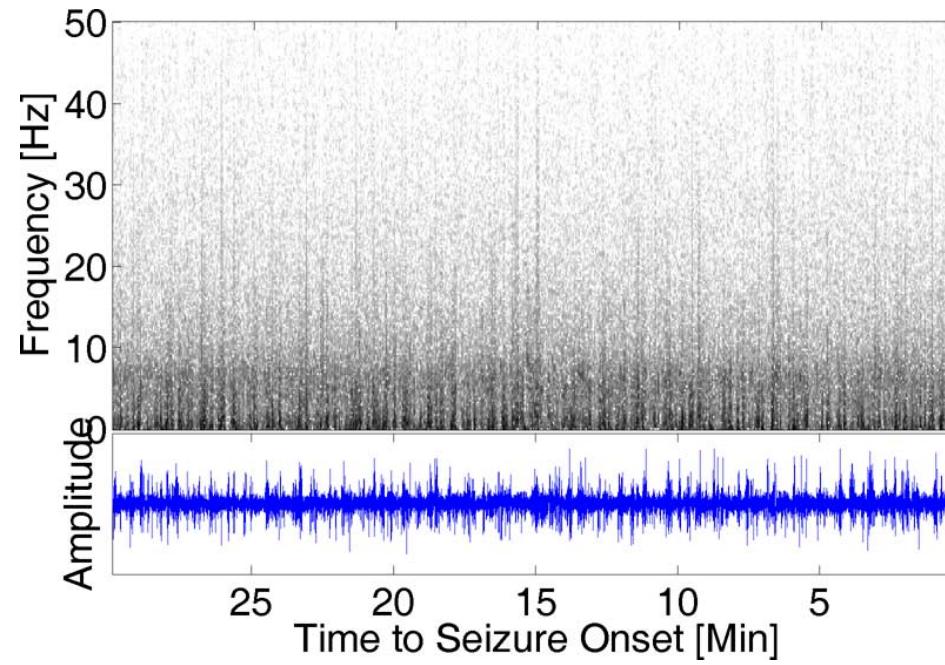
IEEE Transactions on Biomedical Engineering (2015)



Hau-tieng Wu, Ronen Talmon, and Yu-Lun Lo
Assess sleep stage by modern signal processing techniques
IEEE Transactions on Biomedical Engineering (2015)



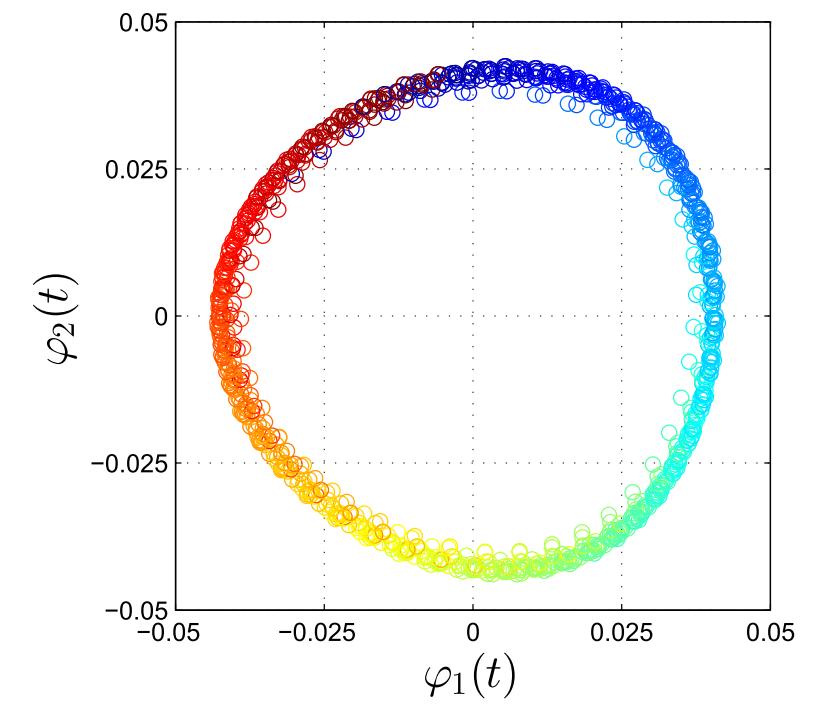
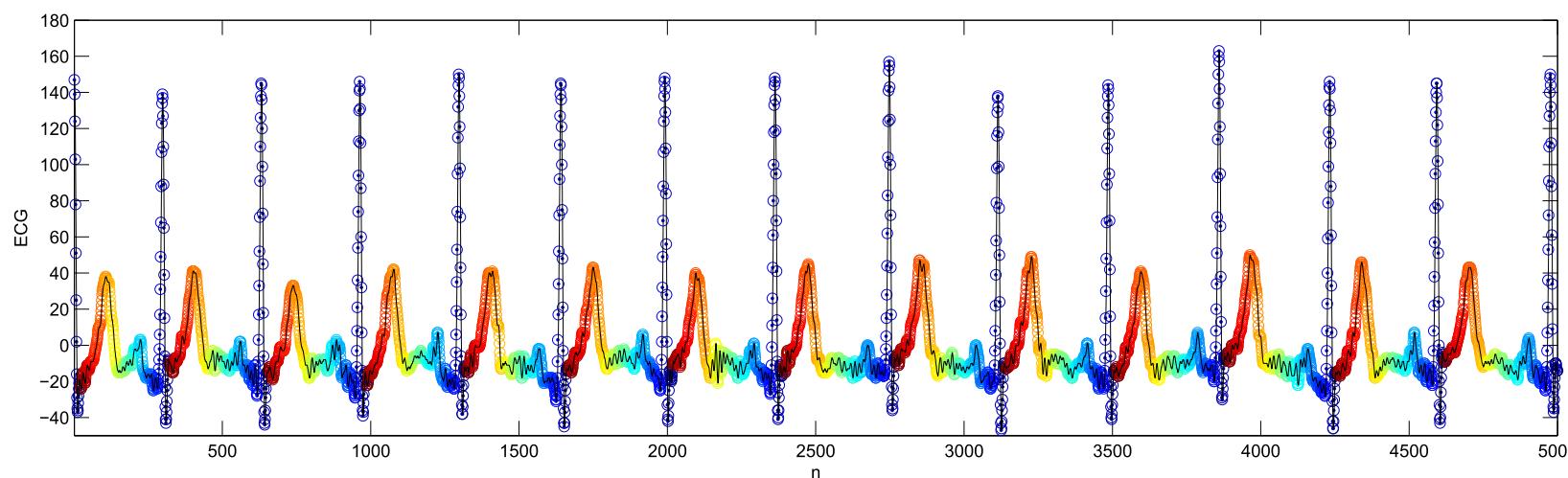
Ronen Talmon, Stéphane Mallat, Hitten Zavari, Ronald R. Coifman
Manifold Learning for Latent Variable Inference in Dynamical Systems
IEEE Transactions in Signal Processing (2015)



Jeremias Sulam, Yaniv Romano, and Ronen Talmon

Dynamical system classification with diffusion embedding for ECG-based person identification

Elsevier Signal Processing (2017)



Concluding remarks

- Diffusion maps and Laplacian eigenmaps are very **similar** techniques
- Tools from multiple **areas**: functional analysis, probability theory, physics, etc.
- Works dealing with large-scale implementation (**quadratic** complexity)
- Questions of **multi-modality** explored recently (Talmon and Wu, 2016)
- Interesting **extensions** to representation learning

Thank you **very much for your attention :)**

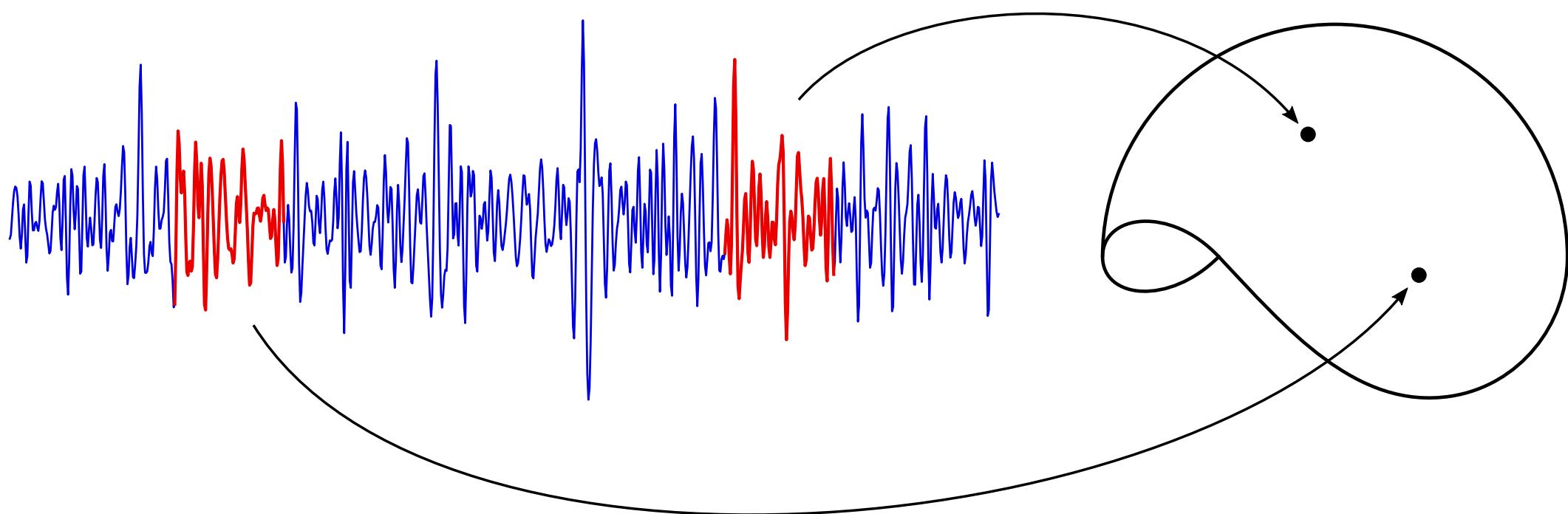
Hope to see some of you again...

CHESS meeting
Monday, 3rd April

Some results on
EEG signal analysis

An introduction to diffusion maps and their use in time series analysis

Pedro Rodrigues – VIBS team meeting, 17th March 2017

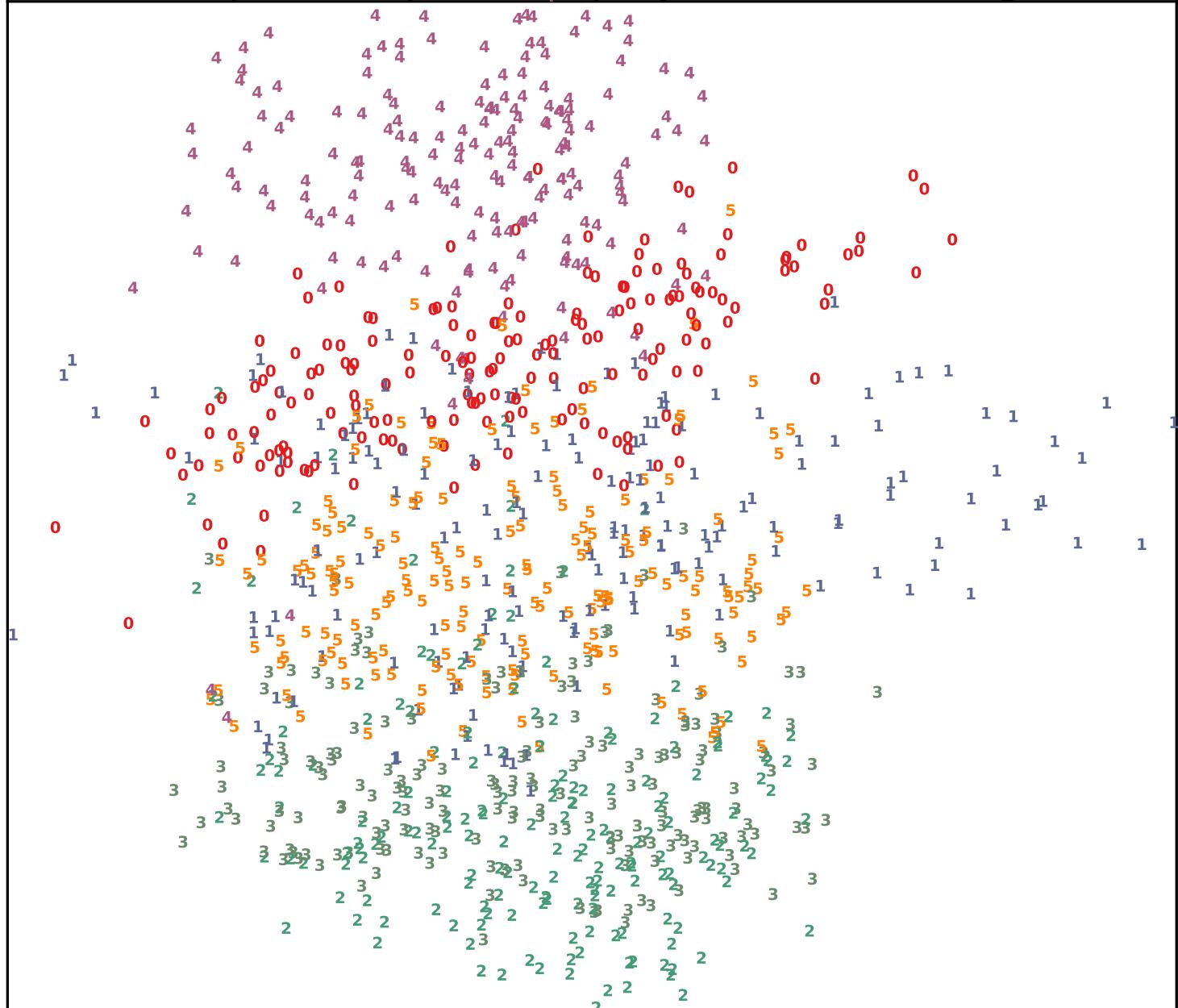


Example with a subset of the MNIST database

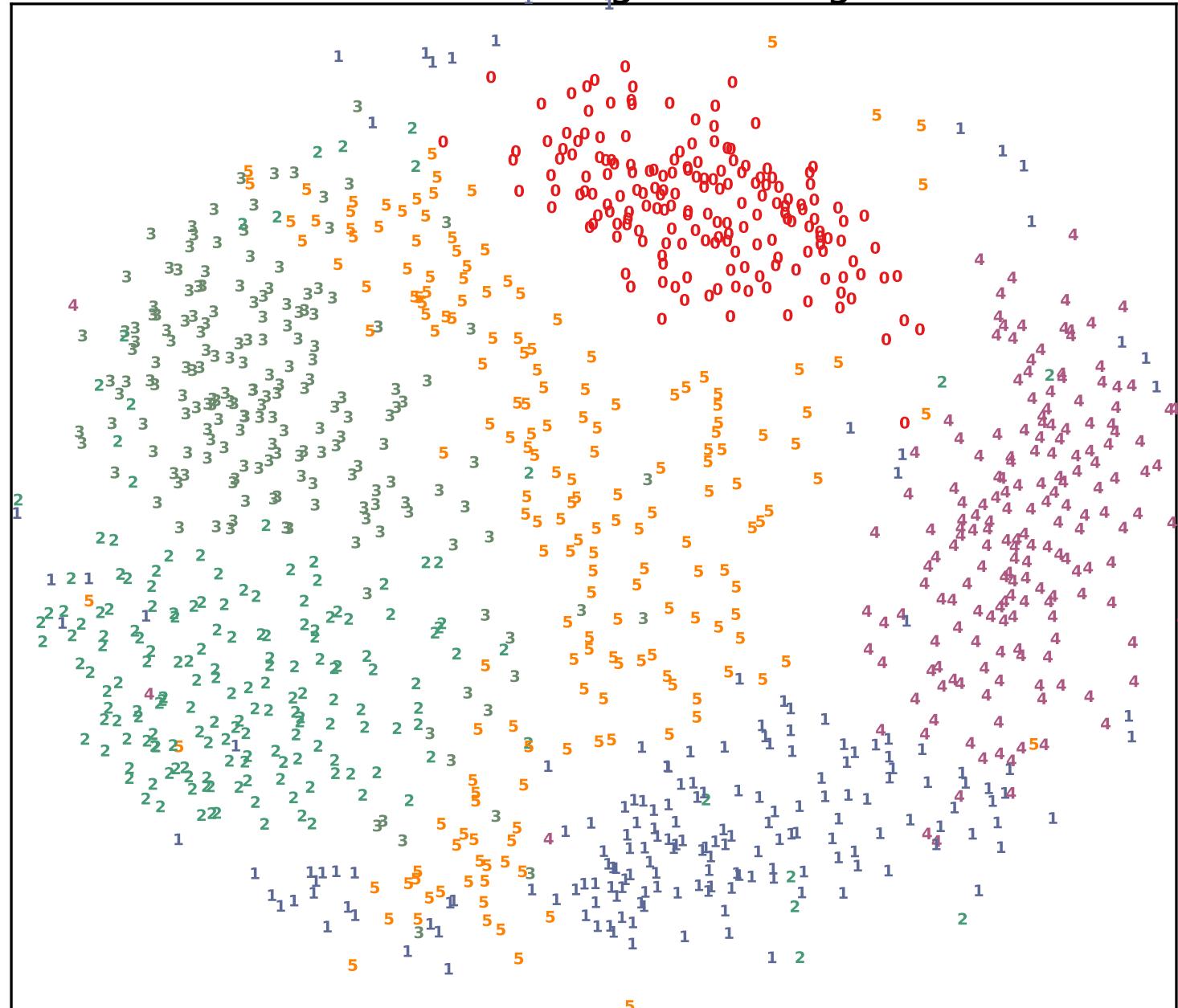
- Handwritten digits
- 400 samples
- 64 x 64 pixel images

0	1	2	3	4	5	0	1	2	3	4	5	0	5
5	5	0	4	1	3	5	1	0	0	2	2	2	0
4	4	1	5	0	5	2	2	0	0	1	3	2	1
3	1	4	0	5	3	1	5	4	4	2	2	5	5
2	3	4	5	0	1	2	3	4	5	0	1	2	3
0	4	1	3	5	1	0	0	2	2	2	0	1	2
1	5	0	5	2	2	0	0	1	3	2	1	3	1
0	5	3	4	5	4	4	2	2	5	5	4	4	0
5	0	4	1	2	3	4	5	0	4	2	3	5	5
3	5	1	0	0	2	2	2	0	1	2	3	3	3
5	2	2	0	0	1	3	2	1	3	1	3	1	4
3	1	5	4	4	2	2	2	5	5	4	4	0	0
0	1	2	3	4	5	0	1	2	3	4	5	0	4
5	1	0	0	1	2	2	0	1	2	3	3	3	4
2	2	0	0	1	3	2	1	4	3	1	4	3	1
1	5	4	4	2	2	2	5	5	4	4	0	3	0
0	1	2	3	4	5	0	1	2	3	4	5	0	4
5	1	0	0	1	2	2	0	1	2	3	3	3	4
2	2	0	0	1	3	2	1	4	3	1	4	3	1
1	5	4	4	2	2	2	5	5	4	4	0	0	1
2	3	4	5	0	1	2	3	4	5	0	5	5	0
0	0	1	2	2	0	1	2	3	3	3	4	4	1
0	0	1	3	2	1	4	3	1	4	3	1	4	0
4	4	2	2	1	5	5	4	4	0	0	1	2	3

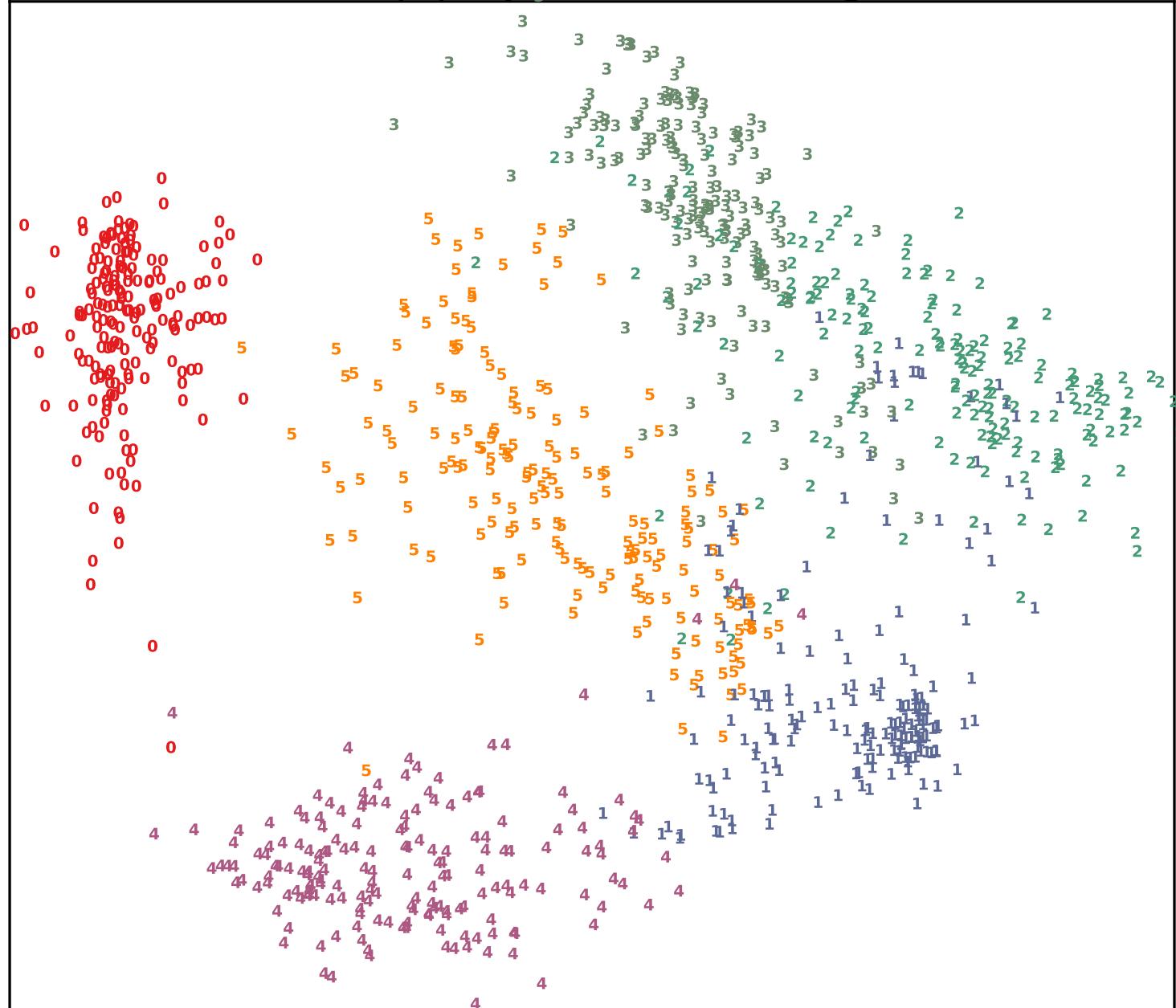
Principal Components projection of the digits



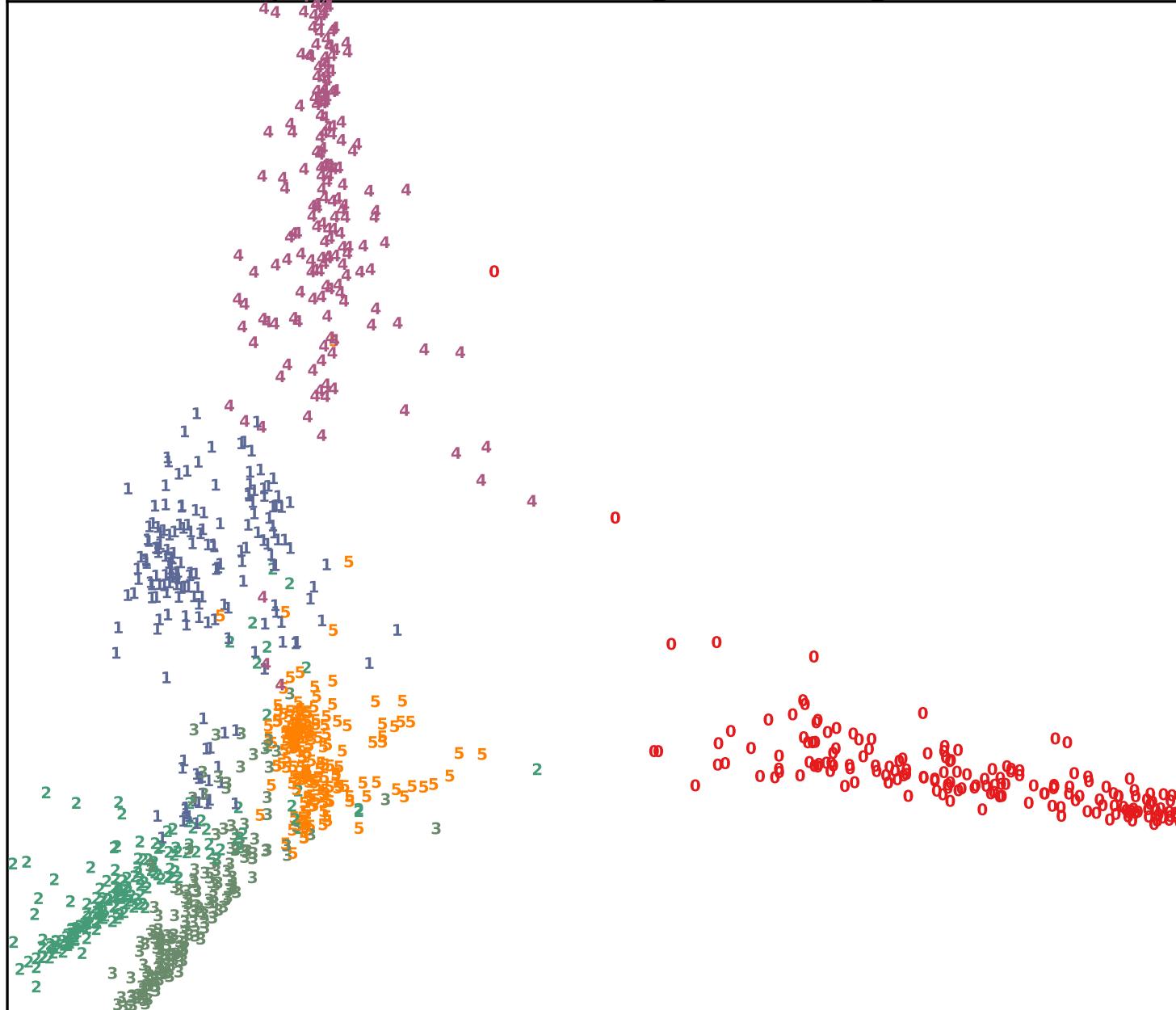
MDS embedding of the digits



Isomap projection of the digits



Spectral embedding of the digits



t-SNE embedding of the digits

