

Parametrization of Datasets with Low-distortion Embeddings

François Meyer

Center for the Study of Brain, Mind and Behavior,
Program in Applied and Computational Mathematics
Princeton University

fmeyer@princeton.edu

<http://www.princeton.edu/~fmeyer>

Random Shapes, Tutorials, IPAM 2007

Acknowledgments

- R.R. Coifman, S. Lafon, M. Maggioni
- M. Ramírez-Vélez, D.S. Barth
- National Institutes of Health
- IPAM MGA program
- IPAM Graduate Summer School: Intelligent Extraction of Information from Graphs and High Dimensional Data

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

4 Classification of EEG recordings

- Explosion of high-dimensional datasets:
web, biology, medicine, etc.
- New tools for data exploration and analysis address issues:
 - ▶ signal of interest: complicated geometry
 - ▶ data corrupted by noise,
 - ▶ algorithm complexity

Example 1: Neuroimaging

- dogma: each brain region is responsible for a specific function
- goal: delineation of functional anatomy in terms of spatial and temporal organization
- method:
 - ▶ very simple cognitive or sensory input stimulus
 - ▶ measure the output signal χ_i at each voxel i inside the brain
 - ▶ detect significant changes in the signal

Example 1: fMRI of Natural Stimuli

- challenge: study the response to complex stimuli (“real life”)
- example: subject watches a movie in the MRI scanner
- discover neuronal networks involved in complex tasks
- how is the analysis performed ?
- size of the problem: 200,000 time series in \mathbb{R}^{1000}

Example 2: Prediction of seizure from EEG

- Electroencephalogram: electrical recordings on the scalp
- seizure \rightarrow time-frequency changes in the signal
- goal: predict the seizure before the onset
- best existing method: brain = nonlinear dynamical system
- neuronal synchrony: fewer independent variables needed to describe the EEG recordings ?
- Is the brain during a seizure a low dimensional dynamical system ?
- size of the problem: 100,000 brain states in \mathbb{R}^{100}

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

4 Classification of EEG recordings

Assumptions about the data

- dataset: large number of internal microscopic variables
→ many degrees of freedom
- at a macroscopic scale: many variables are coupled
→ set of all possible configurations for the signals is low dimensional
- signals varies smoothly as a function of “hidden” variables
→ well defined low dimensional structure

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

4 Classification of EEG recordings

Our definition of the problem

- 1 $T \times N$ dataset $\mathbf{X} = [\mathbf{x}_0 | \mathbf{x}_1 | \cdots | \mathbf{x}_{N-1}]$
- 2 goal: construction of a new parameterization
 $\phi : \mathbb{R}^T \rightarrow \mathbb{R}^K$, with $K \ll T$,
 $\mathbf{x}_i \mapsto \phi(\mathbf{x}_i)$,
- 3 similar signals are mapped to the same region of the atlas:
 - ▶ $\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\| \approx \|\mathbf{x}_i - \mathbf{x}_j\|$

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

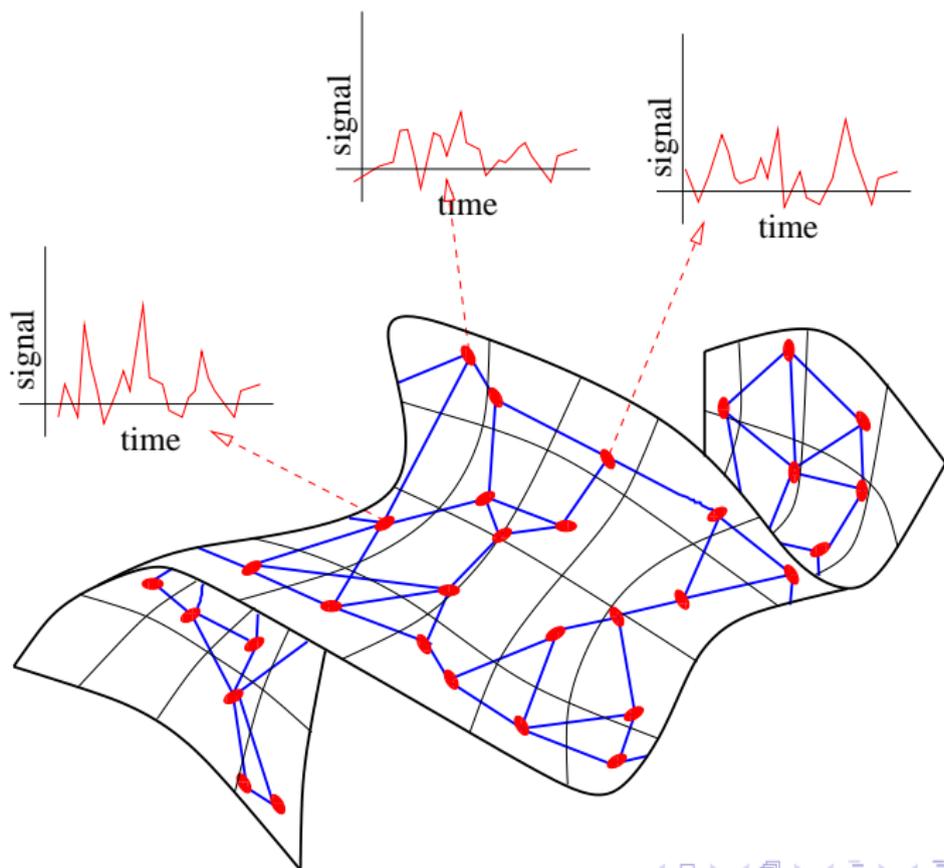
4 Classification of EEG recordings

Construction of the graph

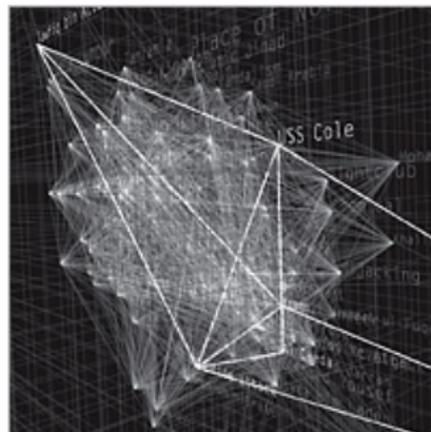
- replace the dataset by a graph G ,
- vertex i of the graph = \mathbf{x}_i
- edges: k nearest neighbors according to
$$\|\mathbf{x}_i - \mathbf{x}_j\| = (\sum_{t=1}^T (x_i(t) - x_j(t))^2)^{1/2}$$
- weight $w_{i,j}$ on the edge $\{i, j\}$: proximity between i and j ,
- for instance,

$$w_{i,j} = \begin{cases} e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2/\sigma^2}, & \text{if } \mathbf{x}_i \text{ is connected to } \mathbf{x}_j, \\ 0 & \text{otherwise.} \end{cases}$$

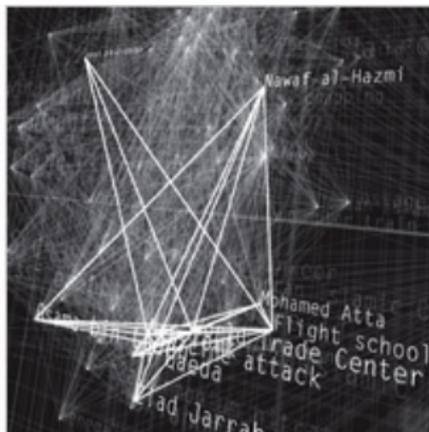
From the dataset to the graph



From the dataset to the graph



OCT. 12, 2000



JUNE 2000 TO JUNE 2001



SEPT. 11, 2001

A random walk on the graph

- weighted graph G , matrix \mathbf{W} , $W_{i,j} = w_{i,j}$
- random walk on the graph with transition probability \mathbf{P} ,

$$P_{i,j} = w_{i,j} / d_i,$$

- $d_i = \sum_j w_{i,j}$: degree of the vertex i , \mathbf{D} diagonal matrix,

$$D_{ii} = d_i = \sum_j W_{i,j}. \quad (1)$$

- $\boldsymbol{\pi} = \frac{1}{\sum_{i,j} w_{i,j}} [d_1, d_2, \dots, d_N]$ stationary distribution

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

4 Classification of EEG recordings

A good distance on the graph

- similarity measure between any two vertices i and j
- distinguish between strongly connected vertices and weakly connected vertices
- solution: average commute time, $\kappa(i, j) = H(j, i) + H(i, j)$
- symmetric version of the average hitting time from i to j ,

$$H(i, j) = E_i[T_j] \quad \text{with} \quad T_j = \min\{n \geq 0; Z_n = j\}.$$

E_i : random walk is started at i

- κ is a distance:
 - 1 $\kappa(i, j) = 0 \implies i = j$,
 - 2 $\kappa(i, j) \leq \kappa(i, k) + \kappa(k, j)$,
 - 3 $\kappa(i, j) = \kappa(j, i)$.

Commute time in Paris



How does $\kappa(i, j)$ compare to $\delta(i, j)$?

- $\kappa(i, j)$ can be compared to the standard distance δ on the graph

Theorem

If i and j are at a distance $\delta(i, j)$ on the graph, then

$$2\delta(i, j) \leq \kappa(i, j) \leq C\delta(i, j),$$

where $C = \max_{i,j} \frac{1}{\pi_i P_{i,j}} = \frac{\sum_{i,j} w_{i,j}}{\min_{i,j} w_{i,j}}$

- Markov chain is reversible, $\pi_i P_{i,j} = \pi_j P_{j,i}$
- C can be large

Your worst commute time in L.A.: Sepulveda Blvd or 405 ?

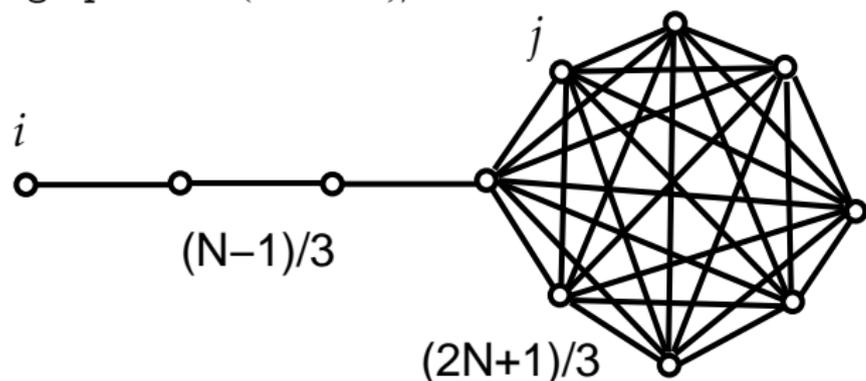


Maximum commute time: lost in the city...

- among all graphs with N vertices,
what is the graph with the largest $\kappa(i, j)$?

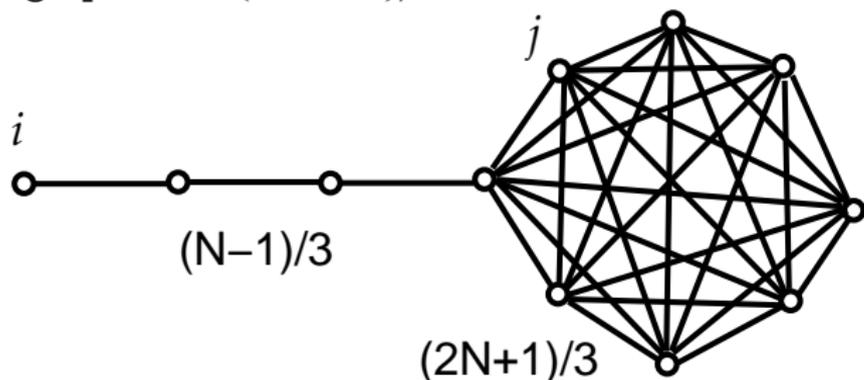
Maximum commute time: lost in the city...

- among all graphs with N vertices, what is the graph with the largest $\kappa(i, j)$?
- lollipop graph: path with $(N - 1)/3$ vertices, complete subgraph with $(2N + 1)/3$ vertices



Maximum commute time: lost in the city...

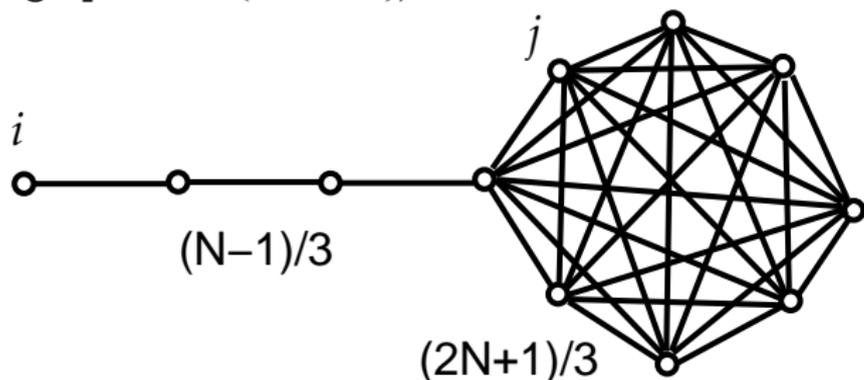
- among all graphs with N vertices, what is the graph with the largest $\kappa(i, j)$?
- lollipop graph: path with $(N - 1)/3$ vertices, complete subgraph with $(2N + 1)/3$ vertices



- $\kappa(i, j) = \frac{4}{27}N^3 + \mathcal{O}(N)$

Maximum commute time: lost in the city...

- among all graphs with N vertices, what is the graph with the largest $\kappa(i, j)$?
- lollipop graph: path with $(N - 1)/3$ vertices, complete subgraph with $(2N + 1)/3$ vertices



- $\kappa(i, j) = \frac{4}{27}N^3 + \mathcal{O}(N)$
- $\delta(i, j) = \frac{2}{3}N$, $C = 2N(2N + 1)/18$
- [Jonasson, 2000]

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

4 Classification of EEG recordings

The spectral connection

- Fundamental matrix $\mathbf{Z} = (I - (\mathbf{P} - \mathbf{\Pi}))^{-1} = I + \sum_{k \geq 1} \mathbf{P}^k - \mathbf{\Pi}$
with $\mathbf{\Pi}^T = [\pi_1 | \cdots | \pi_N]$
- \mathbf{Z} is the Green function of the Laplacian, $I - \mathbf{L}$

Theorem

[Bremaud, 1999] Hitting time $E_i[T_j] = (Z_{j,j} - Z_{i,j})/\pi_j$.

- $E_i[T_j] = 1 + \sum_{k; k \neq j} P_{i,k} E_i[T_k]$
- eigenfunctions ϕ_1, \dots, ϕ_N of

$$\mathbf{D}^{\frac{1}{2}} \mathbf{P} \mathbf{D}^{-\frac{1}{2}}, \quad (2)$$

with eigenvalues $-1 \leq \lambda_N \cdots \leq \lambda_2 < \lambda_1 = 1$.

The spectral connection

- commute time:

$$\kappa(i, j) = \sum_{k=2}^N \frac{1}{1 - \lambda_k} \left(\frac{\Phi_k(i)}{\sqrt{\pi_i}} - \frac{\Phi_k(j)}{\sqrt{\pi_j}} \right)^2. \quad (3)$$

- define an embedding

$$i \mapsto I_k(i) = \frac{1}{\sqrt{1 - \lambda_k}} \frac{\Phi_k(i)}{\sqrt{\pi_i}}, \quad k = 2, \dots, N \quad (4)$$

- Euclidean distance on the image of the embedding
= commute time

$$\kappa(i, j) = \|I(i) - I(j)\|^2 = \sum_{k=2}^N |I_k(i) - I_k(j)|^2$$

1 Introduction

- Assumptions about the data
- Our definition of the problem

2 Constructing a new parametrization

- A random walk on the dataset
- A new way to measure distances

3 The spectral connection

- Spectral graph theory
- From commute time to spectral geometry

4 Classification of EEG recordings

The Laplacian connection

- Φ_k is also an eigenfunction of the Laplacian

$$\mathcal{L} = I - \mathbf{D}^{\frac{1}{2}} \mathbf{P} \mathbf{D}^{-\frac{1}{2}},$$

with the eigenvalues $\beta_k = 1 - \lambda_k$.

- Φ_k minimizes the “distortion”

$$\min_{\Phi, \|\Phi\|=1} \frac{\sum_{[i,j]} w_{i,j} (\Phi(i) - \Phi(j))^2}{\sum_i d_i \Phi^2(i)}$$

with Φ_k orthogonal to $\{\Phi_0, \Phi_1, \dots, \Phi_{k-1}\}$.

- Laplacian eigenmaps [Belkin and Niyogi, 2003]

The diffusion distance connection

- [Lafon, 2004, Coifman and Lafon, 2006]
- diffusion distance,

$$D_t^2(i, j) = \sum_{k=2}^N \lambda_k^{2t} \left(\frac{\Phi_k(i)}{\sqrt{\pi_i}} - \frac{\Phi_k(j)}{\sqrt{\pi_j}} \right)^2. \quad (5)$$

- commute time = sum of the diffusion distance at all scale t

$$\sum_{t=0}^{\infty} D_{t/2}^2(i, j) = \sum_{k=2}^N \frac{1}{1 - \lambda_k} \left(\frac{\Phi_k(i)}{\sqrt{\pi_i}} - \frac{\Phi_k(j)}{\sqrt{\pi_j}} \right)^2 = \kappa(i, j)$$

The spectral geometry connection

- data sampled on a n -dimensional manifold \mathcal{M}
- \mathcal{M} embedded by its heat kernel $K_{\mathcal{M}}(t, x, y)$

Theorem

[Bérard et al., 1994]

$$\psi_t : \mathcal{M} \mapsto l^2(\mathbb{R}) \quad (6)$$

$$x \mapsto \left\{ \sqrt{2}(4\pi)^{n/4} t^{(n+2)/4} e^{-\lambda_j t/2} \phi_k(x) \right\}_{k \geq 1} \quad (7)$$

$\forall t > 0$, the map ψ_t is an embedding of \mathcal{M} into $l^2(\mathbb{R})$.

- scale parameter t : similar to diffusion distance

The spectral geometry connection

ψ_t : composition of

- 1 embedding of \mathcal{M} by the heat kernel:
each point on \mathcal{M} is mapped to a bump function.

$$\mathcal{M} \mapsto L^2(\mathcal{M}) \quad (8)$$

$$x \mapsto K_{\mathcal{M}}(t/2, x, \cdot) \quad (9)$$

- 2 isometry given by the choice of basis, $\{\Phi_1, \Phi_2, \dots\}$, of $L^2(\mathcal{M})$,
each function of $L^2(\mathcal{M})$ is expanded into the basis of
eigenfunctions of the Laplace-Beltrami operator

$$L^2(\mathcal{M}) \mapsto l^2(\mathbb{R}) \quad (10)$$

$$f \mapsto \{\langle f, \Phi_k \rangle\}_{k \geq 1} \quad (11)$$

Algorithm 1: Construction of the embedding

Input:

- ▶ $\mathbf{x}_i(t)$, $t = 0, \dots, T - 1$, $i = 1, \dots, N$,
- ▶ σ ; n_n number of nearest neighbors.
- ▶ K : number of eigenfunctions.

Algorithm:

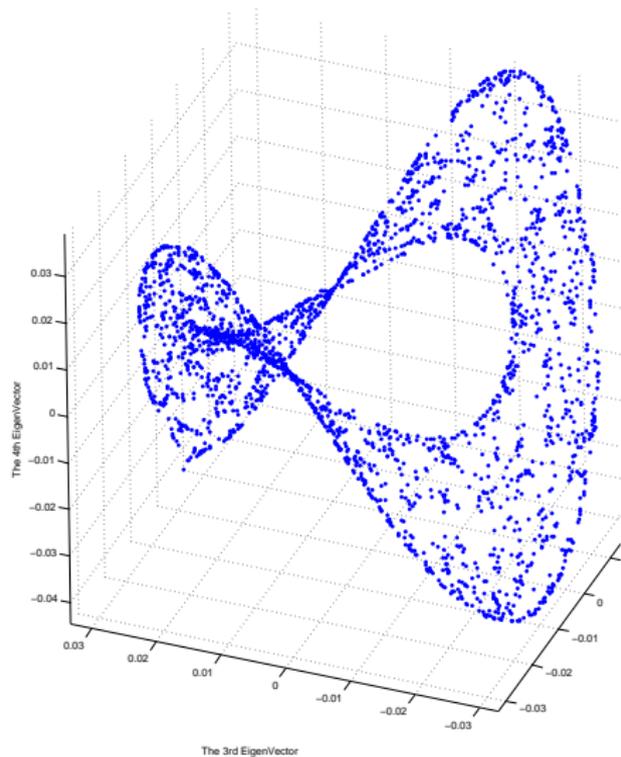
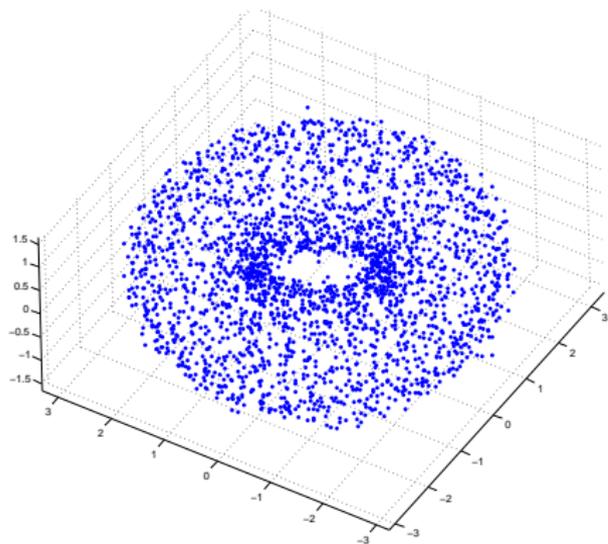
- 1 construct the graph defined by the n_n nearest (according to $\|\mathbf{x}_i - \mathbf{x}_j\|$) neighbors of each \mathbf{x}_i
- 2 compute \mathbf{P}
- 3 find the first K eigenfunctions, Φ_k , of $\mathbf{D}^{\frac{1}{2}} \mathbf{P} \mathbf{D}^{-\frac{1}{2}}$

Output:

 For all \mathbf{x}_i

- ▶ new co-ordinates of \mathbf{x}_i : $\left\{ \frac{1}{1-\lambda_k} \frac{\Phi_k(i)}{\sqrt{\pi_i}} \right\}$, $k = 2, \dots, N$

A toy example



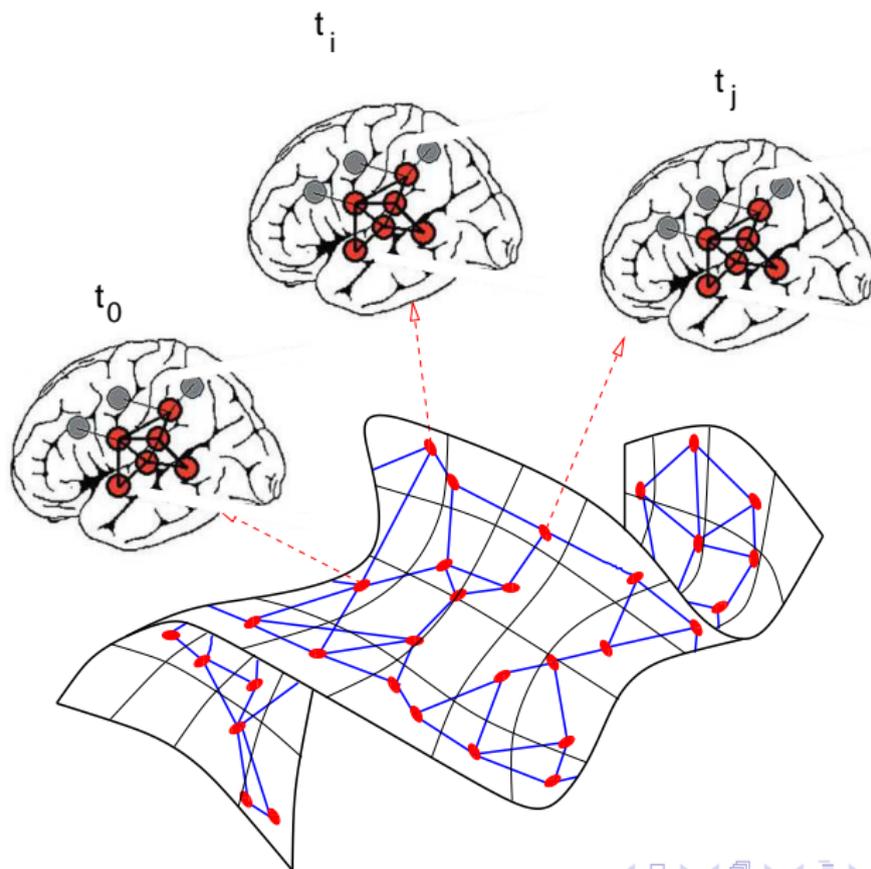
Classification of EEG recordings

- classification of EEG recordings into baseline and ictal states
- Hypothesis: we can find a lower dimensional representation for classification
- More details can be found in
[Ramírez-Vélez et al., 2006, Meyer and Shen, 2007]

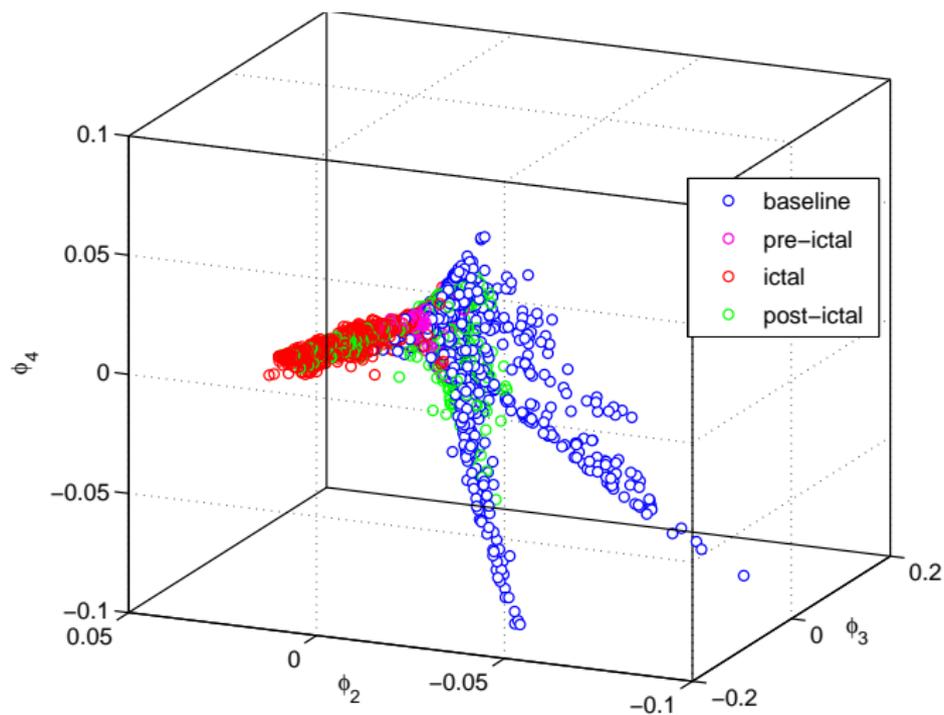
Human dataset

- scalp electroencephalograms
- 55 electrode channels, lowpass filtered at 256 Hz.
- baseline, pre-ictal, ictal, and post-ictal time segments
- each node of the graph is in \mathbb{R}^{55}

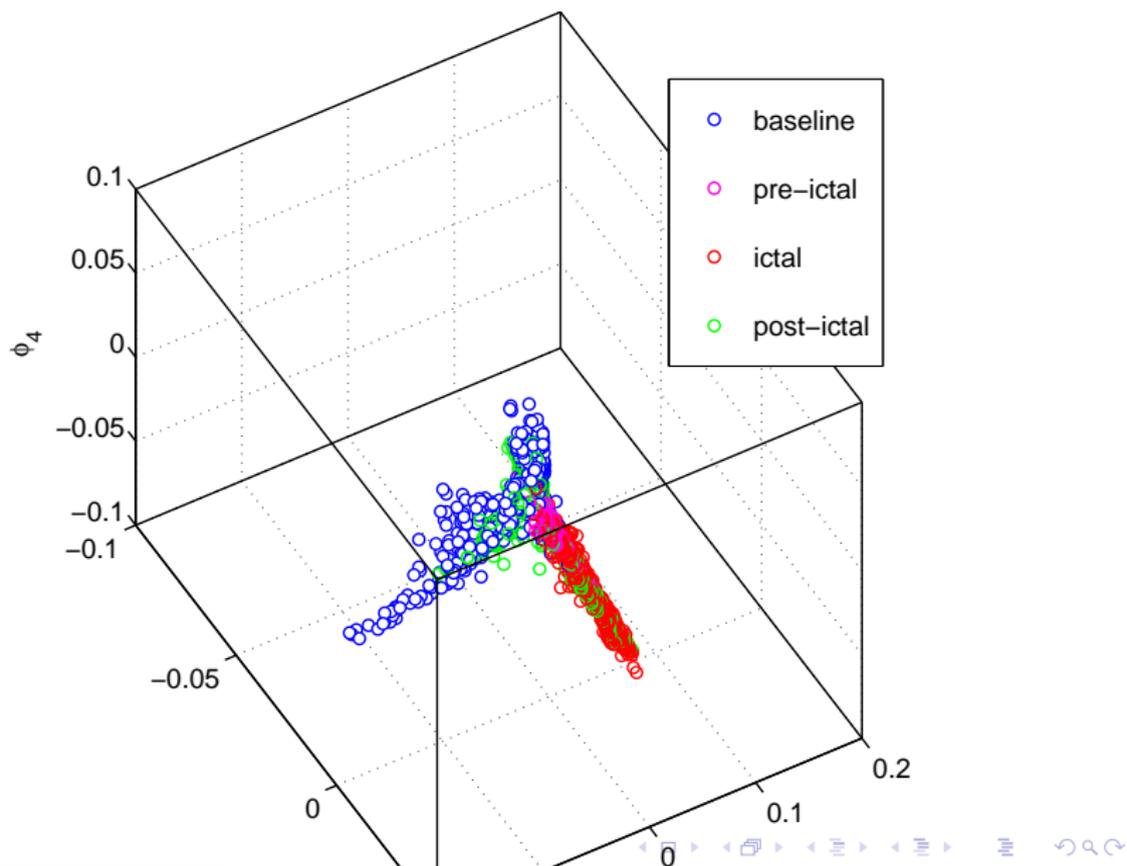
The graph of the brain dynamics



Embedding



Embedding (2)



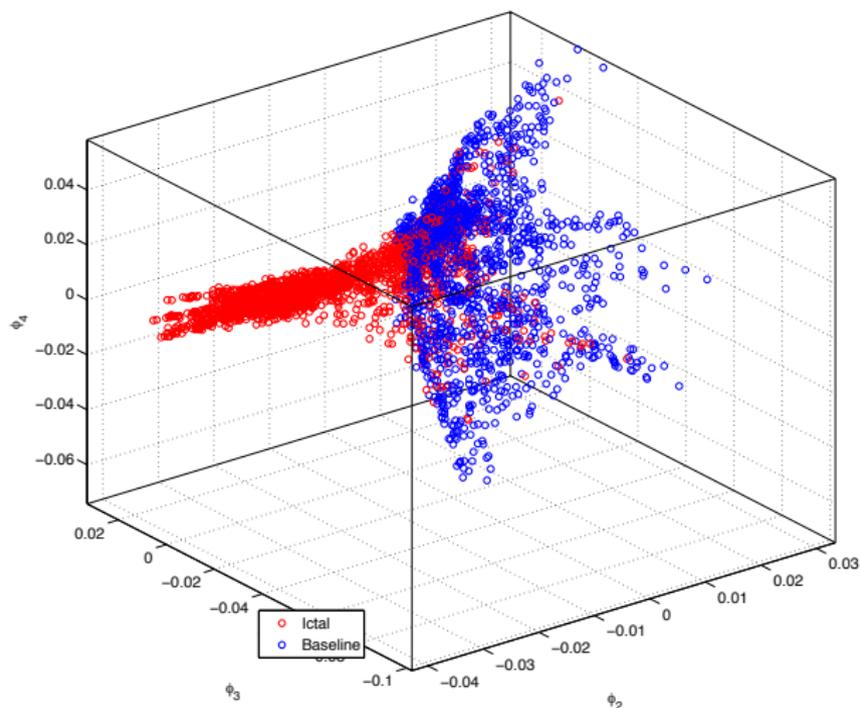
Classification

- $d = 10$ dimensions
- 10-fold cross validation

Table: Classification error (kernel ridge regression)

	Baseline	Ictal	Total
Raw Data	93.20	70.40	81.80
PCA	81.60	78.80	82.20
Random walk	100	82.80	91.40

Classification



Estimated labels: red=ictal, blue=baseline

Curse of dimensionality:

- Fast nearest neighbors in high dimension
- Eigensolvers for large ($N = 10^5 - 10^6$) sparse matrices

Open questions

- much faster eigensolvers are needed...Matlab blows up for $N > 10,000$
- real time update ϕ_k with new incoming data ?
- ϕ_k : sensitive to σ and n_n
- how many new co-ordinates (local dimension) ?



Belkin, M. and Niyogi, P. (2003).

Laplacian eigenmaps for dimensionality reduction and data representation.

Neural Computations, 15:1373–1396.



Bérard, P., Besson, G., and Gallot, S. (1994).

Embeddings Riemannian manifolds by their heat kernel.

Geometric and Functional Analysis, 4(4):373–398.



Bremaud, P. (1999).

Markov Chains.

Springer Verlag.



Coifman, R. and Lafon, S. (2006).

Diffusion maps.

Applied and Computational Harmonic Analysis, 21:5–30.



Jonasson, J. (2000).

Lollipop graphs are extremal for commute times.

Random Structures and Algorithms, 16(2):131–142.

 Lafon, S. (2004).
Diffusion maps and geometric harmonics.
PhD thesis, Yale University, New Haven.

 Meyer, F. and Shen, X. (2007).
Exploration of high dimensional biomedical datasets with
low-distortion embedding.
*In Proc. Data Mining for Biomedical Informatics Workshop,
7th SIAM International Conference on Data Mining.*
To appear.

 Ramírez-Vélez, M., Staba, R., Barth, D., and Meyer, F. (2006).
Nonlinear classification of EEG data for seizure detection.
*In Proc. IEEE International Symposium on Biomedical
Imaging: Macro to Nano*, pages 956–959.