

NONLINEAR EMBEDDINGS

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

LOWER DIMENSIONAL (METRIC) EMBEDDINGS

Spectral connectivity analysis (SCA) is a general process for finding lower dimensional structure in the data

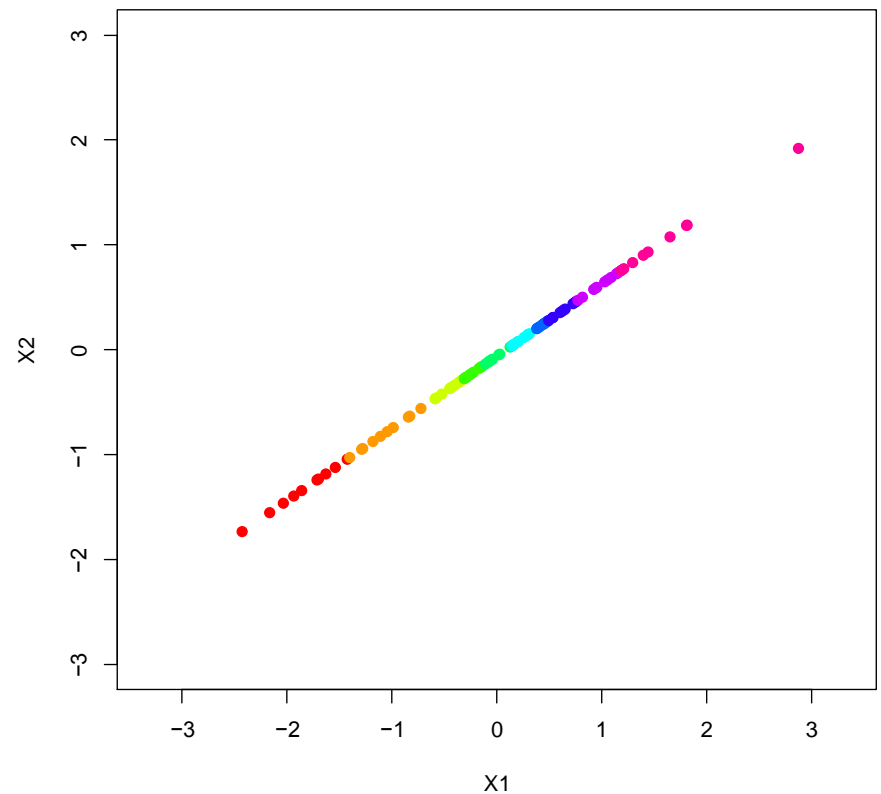
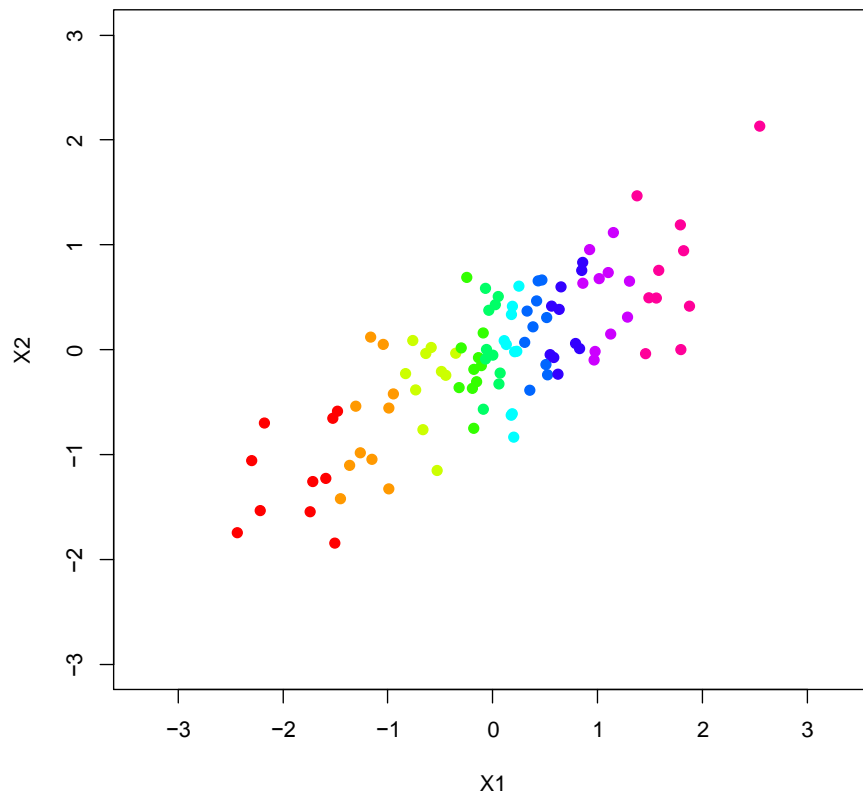
It can be...

- Linear or nonlinear
- Used for dimension reduction or feature creation
- PCA, Fisher discriminant analysis, Locally linear embeddings, Hessian eigenmaps, **Laplacian eigenmaps**, **kernel PCA**
- Useful as an input to classification, clustering, and regression approaches

Let's take one last look at PCA before proceeding

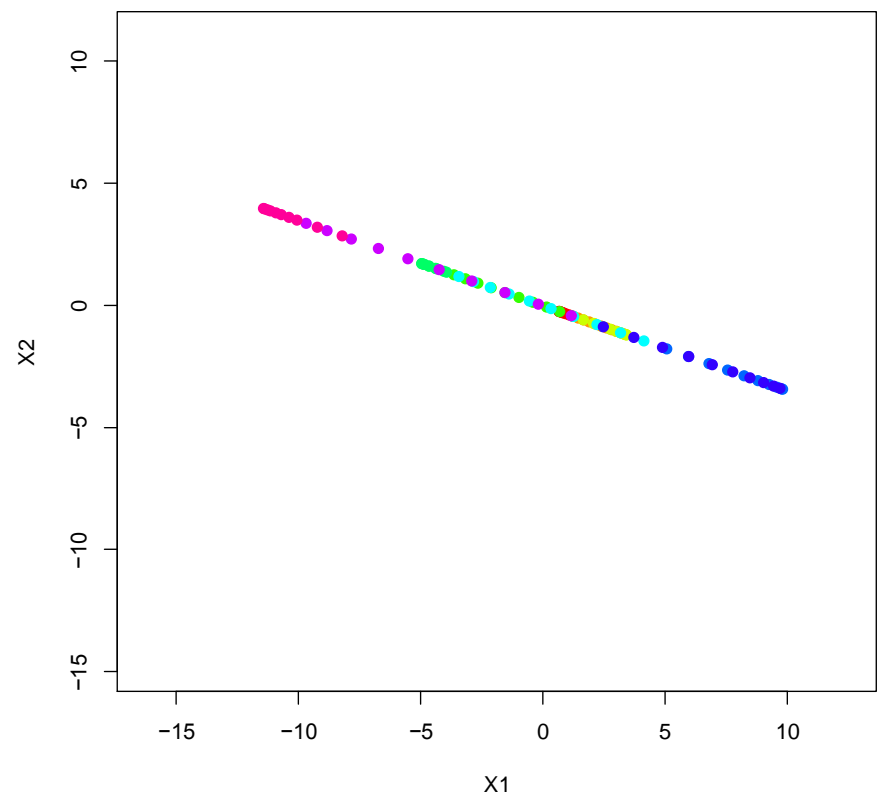
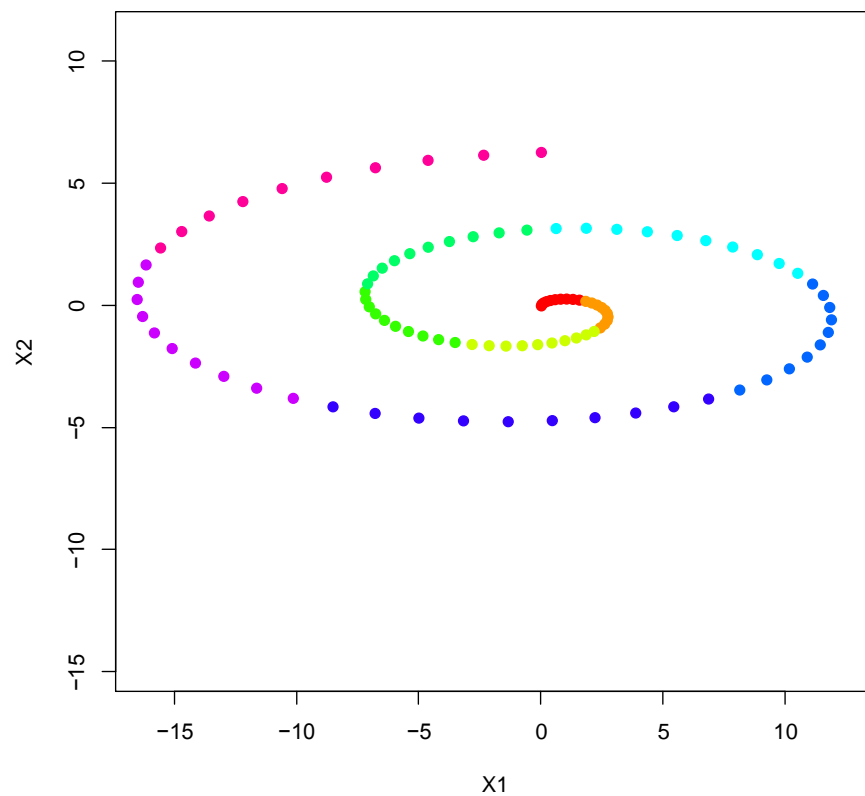
WHEN PCA WORKS WELL

PCA can do effective dimension reduction (that is, explain most of the data with $m < p$ components) as long as the data can be efficiently represented as 'lines' (or planes, or hyperplanes). So, in two dimensions:



WHEN PCA DOESN'T WORK WELL

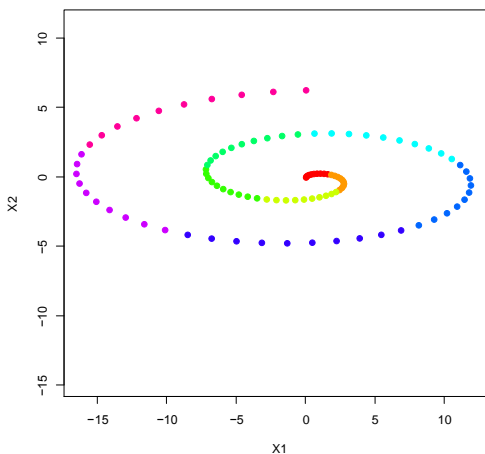
What about other data structures? Again in two dimensions



Here, we have failed miserably.

EXPLANATION

- PCA wants to minimize distances (equivalently maximize variance). This means it **slices** through the data at the **meatiest** point, and then the next one, and so on. If the data are 'curved' this is going to induce artifacts.
- PCA also looks at things as being **close** if they are near each other in a Euclidean sense
[this is essentially all covariance is].
- On the spiral, our intuition says that things are 'close' only if the distance is constrained to go along the curve. In other words, purple and blue are close, blue and red are not.



PCA AND COVARIANCE

PCA: Find the directions of greatest variance. This doesn't on its face seem like it maintains correlations, but observe:

$$\text{var}([a, b]^T X) = \text{var}(ax_1 + bx_2) = a^2 \text{Var}(x_1) + b^2 \text{Var}(x_2) + 2ab \text{Cov}(x_1, x_2)$$

If we standardize the matrix, then this reduces to $\text{MAX}_{\alpha} \alpha^T X$

$$\text{var}(ax_1 + bx_2) = a^2 + b^2 + 2ab \text{Cov}(x_1, x_2) \quad \text{s.t.} \quad \|\alpha\|_2^2 = 1$$

This gets maximized over $a^2 + b^2 = 1$.

- If $\text{Cov}(x_1, x_2) \approx 0$, then this gets maximized by any $a^2 + b^2 = 1$ (it doesn't matter)
- If $\text{Cov}(x_1, x_2) \approx 1$, then this gets maximized by setting $a = b = 1/\sqrt{2}$

So, in either case, we are really maintaining **correlations**

Correlation is fundamentally a linear phenomenon

GRAPHICAL EXAMPLE OF THE PHENOMENON

```

library(mvtnorm)
sigma = matrix(c(1,sig,sig,1),nrow=2)
nsweep = 1000
outcome = matrix(0,nrow=nsweep,ncol=2)
for(sweep in 1:nsweep){
  x = rmvnorm(200,c(0,0),sigma)
  out.pca = prcomp(x,center=T,scale=F)
  outcome[sweep,] = out.pca$rotation[,1]
}
plot(outcome,xlab='PC1',ylab='PC2')

```

$$\text{SIG} = 0$$

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\text{SIG} \sim 0.999$$

$$\Sigma = \begin{bmatrix} 1 & .999 \\ .999 & 1 \end{bmatrix}$$

"V"

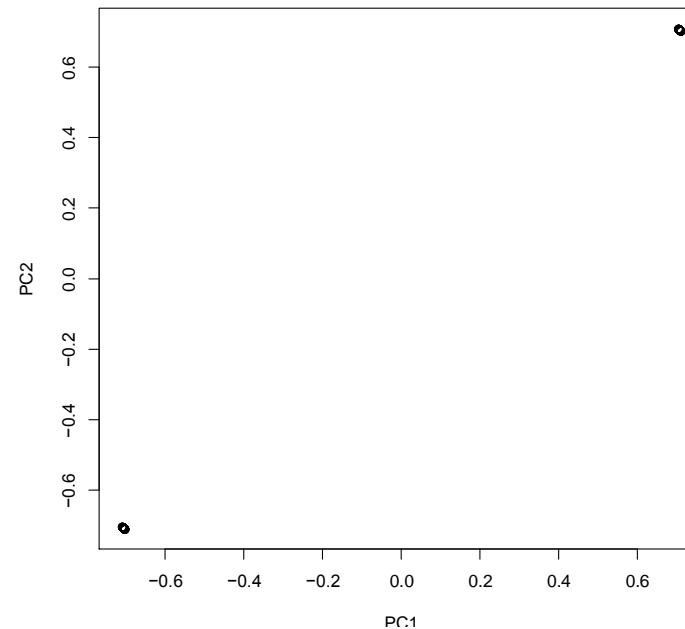
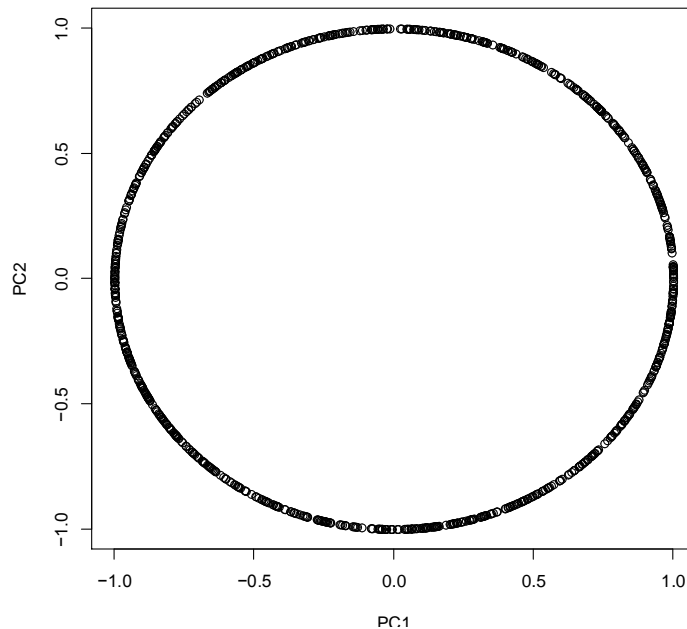


FIGURE: Left: $\text{sig} = 0$. Right: $\text{sig} = .999$

Nonlinear embeddings

KERNEL PCA (KPCA)

$$\tilde{\Sigma} \propto \tilde{\mathbf{X}}^T \hat{\mathbf{X}} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^T (\mathbf{I} - \mathbf{M}) \mathbf{X} \quad \Sigma_1 = (U \% * \% D) [1, 1] \in \mathbb{R}^n$$

Classical PCA comes from $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{M}\mathbf{X} = \underbrace{\mathbf{U}}_{\substack{\text{SCORES} \\ \gamma_{11} \\ \gamma_{21} \\ \gamma_{31}}} \underbrace{\mathbf{D}}_{\substack{\text{LOADINGS} \\ \gamma_{11} \\ \gamma_{21} \\ \gamma_{31}}} \mathbf{V}^T$, where

$\mathbf{M} = \mathbf{1}\mathbf{1}^T / n$ and $\mathbf{1} = (1, 1, \dots, 1)^T$

However, we can just as easily get it from the **outer product**

$$\mathbb{K} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = (\mathbf{I} - \mathbf{M}) \underbrace{\mathbf{X}\mathbf{X}^T}_{\substack{\rightarrow \\ [\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle]_{1 \leq i, i' \leq n}}} (\mathbf{I} - \mathbf{M}) = \mathbf{U}\mathbf{D}^2\mathbf{U}^T \in \mathbb{R}^{n \times n}$$

The intuition behind KPCA is that \mathbb{K} is an expansion into a kernel space, where

$$\mathbb{K}_{i,i'} = k(\tilde{\mathbf{X}}_i, \tilde{\mathbf{X}}_{i'}) = \langle \tilde{\mathbf{X}}_i, \tilde{\mathbf{X}}_{i'} \rangle$$

REMINDER: Anytime we see an inner product, we can kernelize it

KERNEL PCA

Following this intuition, the approach is simple:

1. Specify a kernel k
(e.g. $k(X, X') = \exp\{-\gamma^{-1} \|X - X'\|_2^2\}$)
2. Form $K_{i,j'} = k(X_i, X_{i'})$
3. Standardize and get eigenvector decomposition

$$\mathbb{K} = (I - M)K(I - M) = UD^2U^\top$$

This implicitly finds the inner product:

$$k(X_i, X_{i'}) = \langle \phi(X_i), \phi(X_{i'}) \rangle$$

However, we need only specify the **kernel**

KERNEL PCA

The scores are still $Z = UD$

The q^{th} KPCA score is (up to centering)

$$Z_{iq} = \sum_{i'=1}^n \beta_{i'q} k(X_i, X_{i'})$$

$$K = U D^2 U^T$$

where $\beta_{i',q} = u_{i'q}/d_q$

NOTE: As we don't explicitly generate the feature map, there are no loadings

Small detour

REPRODUCING KERNEL HILBERT SPACE

REMINDER: Mercer's theorem assures us that

$$k(X, X') = \sum_{j=1}^{\infty} \theta_j \phi_j(X) \phi_j(X')$$

sym matrix

$$\Phi \Theta \Phi^T$$

$$\langle \phi_i, \phi_k \rangle = \delta_{jk} = \sum_{j=1}^{\infty} \theta_j \psi_j \psi_j^T$$

Here, the system $(\phi_j)_{j=1}^{\infty}$ spans a space \mathcal{H}_k

The function space \mathcal{H}_k is known as a **reproducing kernel Hilbert space (RKHS)**

It can also be thought of as roughly

$$\mathcal{H}_k = \left\{ f : f(X) = \sum_{i=1}^n \beta_i k(X, X_i) \right\}$$

Which has a special inner product

$$\langle f, f \rangle_{\mathcal{H}_k} = \|f\|_{\mathcal{H}_k}^2 = \sum_{j=1}^{\infty} f_j^2 / \theta_j < \infty$$

$f_j = \langle f, \phi_j \rangle$

REPRODUCING KERNEL HILBERT SPACE

Writing

$$f(X) = \sum_{i=1}^n \beta_i k(X, X_i)$$

The terms $k(X, X_i)$ are the **representers**, as

$$\langle k(\cdot, X), f \rangle_{\mathcal{H}_k} = f(X)$$

and \mathcal{H}_k is called a **reproducing kernel Hilbert space** (RKHS) as

$$\langle k(\cdot, X), k(\cdot, X') \rangle_{\mathcal{H}_k} = k(X, X')$$

NOTE: For kernel methods, we are generalizing the finite dimensional Euclidean inner product

$$\langle X, X' \rangle = X^\top X'$$

KERNEL METHODS VIA REGULARIZATION

After specifying a kernel function k , we can define an estimator via

\mathcal{H}_k

$$\hat{\mathbb{P}} \ell_f \stackrel{\text{SQUARED ERROR}}{=} \frac{1}{n} \sum (Y_i - f(X_i))^2 \quad \min_{f \in \mathcal{H}_k} (\hat{\mathbb{P}} \ell_f) + \lambda \|f\|_{\mathcal{H}_k}^2$$

This is a (potentially) infinite dimensional optimization problem
(hard, especially with a computer)

It can be shown that the solution has the form

$$\hat{f}(X) = \sum_{i=1}^n \beta_i k(X, X_i)$$

(This is known as the **representer theorem**)

Back to KPCA

KERNEL PCA

REMINDER: To get the first PC in classical PCA, we want to solve

$$\max_{\alpha} \mathbb{V} \alpha^{\top} X \quad \text{subject to} \quad \|\alpha\|_2^2 = 1$$

Translate this into the kernel setting, and we are trying to solve

$$\max_{g \in \mathcal{H}_k} \mathbb{V} g(X) \quad \text{subject to} \quad \|g\|_{\mathcal{H}_k} = 1$$

The **representer theorem** states that a solution to this problem is

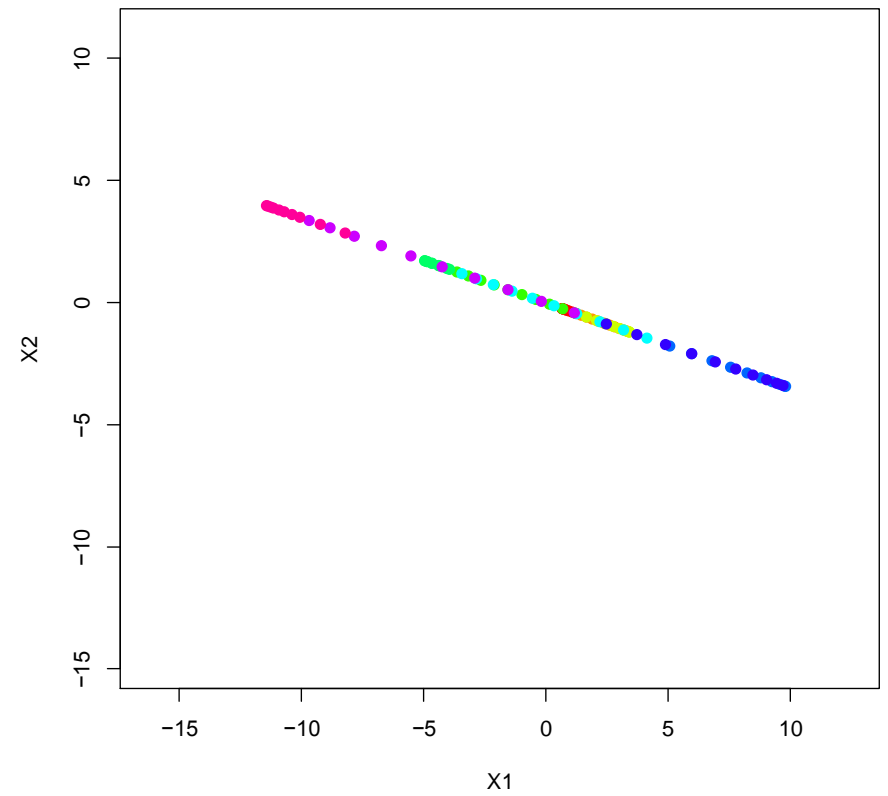
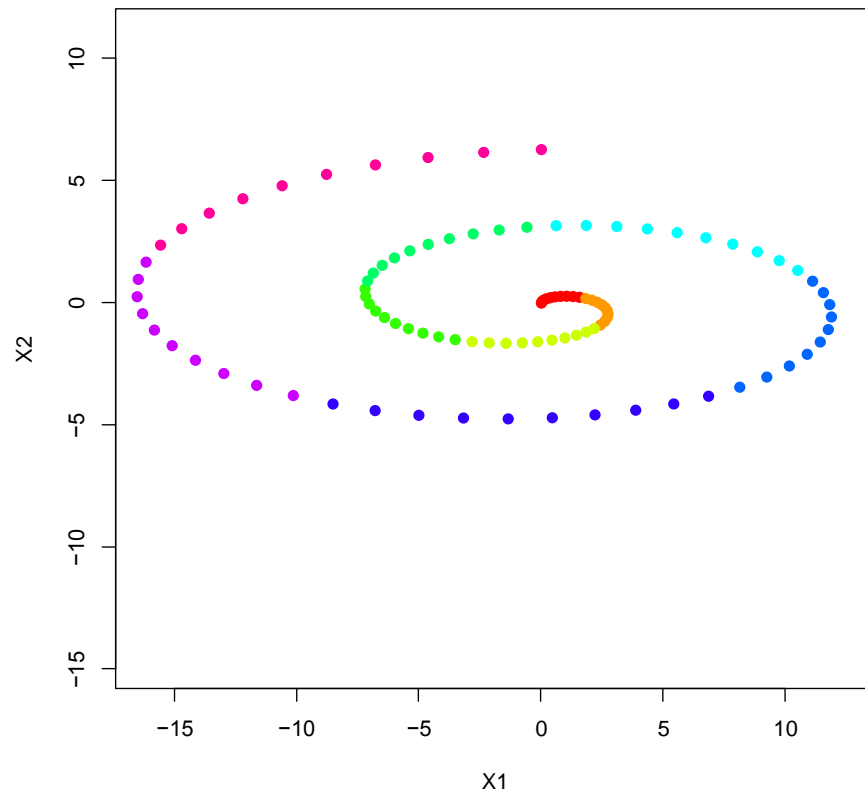
$$\hat{g}(X) = \sum_{i=1}^n \beta_i k(X, X_i) \in \beta^{\top} K$$

Compare

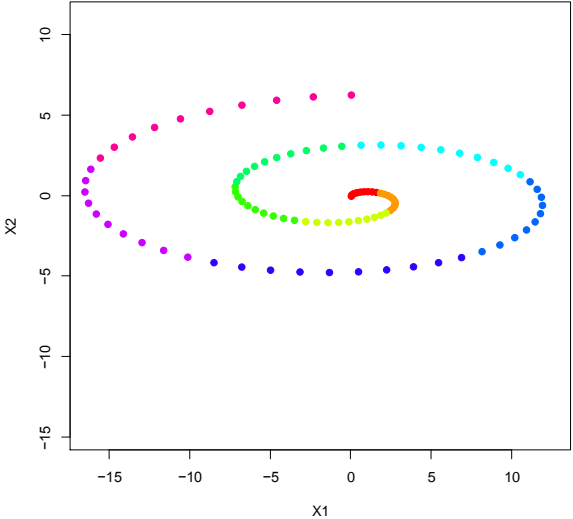
$$Z_{iq} = \sum_{i'=1}^n \beta_{i'q} k(X_i, X_{i'})$$

where $\beta_{i',q} = u_{i'q} / d_q$

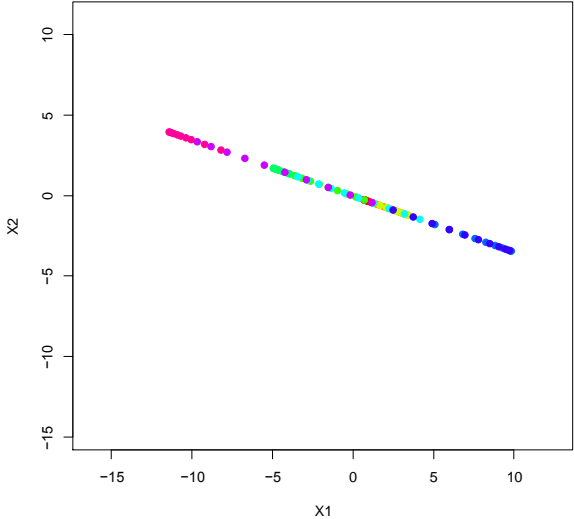
KPCA: SOME RESULTS



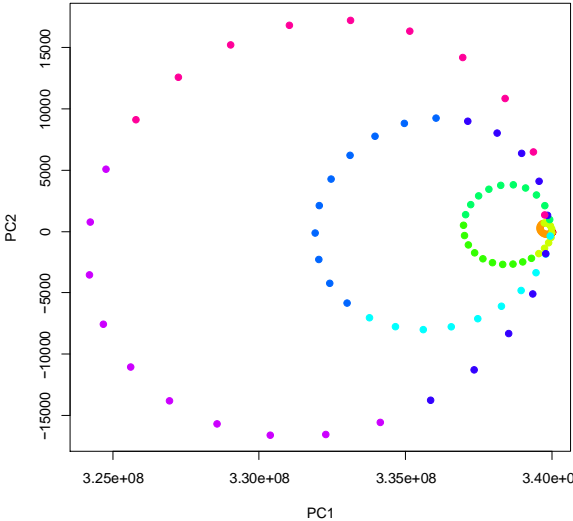
KPCA: SOME RESULTS



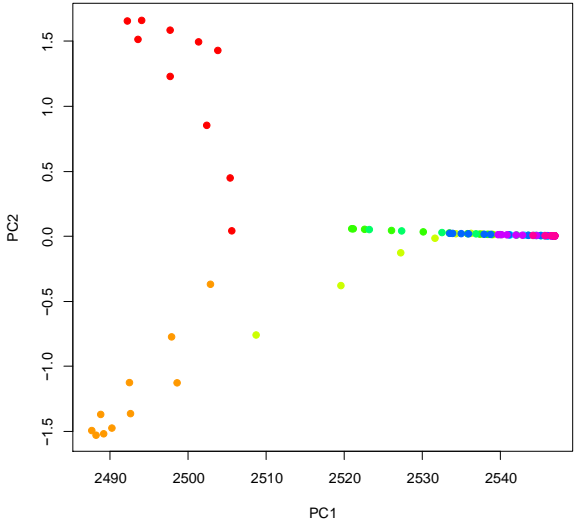
Data



PCA



KPCA: polynomial(2)



KPCA: gaussian(10)

SEMISUPERVISED LEARNING IN PRACTICE

Looking at:

$$Z_{iq} = \sum_{i'=1}^n \beta_{i'q} k(X_i, X_{i'})$$

this is only defined at our observed features

Write

65" 68" 74"

- $\mathcal{D}_{train} = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$
- $\mathcal{D}_{test} = \{(X_1^*, Y_1^*), \dots, (X_{n^*}^*, Y_{n^*}^*)\}$

Two common scenarios are

1. We are given \mathcal{D}_{train} and $X_1^*, \dots, X_{n^*}^*$ to build \hat{f}
2. We are given only \mathcal{D}_{train} to build \hat{f}

CASE 1

We are given \mathcal{D}_{train} and $X_1^*, \dots, X_{n^*}^*$ to build \hat{f}

Then we can just use straight forward KPCA

(Or any unsupervised learning step)

1. Form \mathbb{K} on \mathcal{D}_{train} and $X_1^*, \dots, X_{n^*}^*$
2. Get UD
3. Pass $Z_q = UD[:, 1 : q]$ to train \hat{f}
4. Get $\hat{Y} = \hat{f}(Z_q) \in \mathbb{R}^{\Lambda + \Lambda^*}$

CASE 2

We are given only \mathcal{D}_{train} to build \hat{f}

Now, we don't know the **coordinates** of X_1^*, \dots, X_{n^*} in the representation space

$$\mathbb{K} = U D^2 U^T = (U D)(D U^T)$$

To get a new observation X^* embedded into this representation:

$$Z_0^T = [z_{01}, \dots, z_{0n}] \quad Z_0 = D^{-1} U^T (I - M) [k^* - K \mathbf{1}/n]$$

$$[UD][1, \dots, n] = [z_{01}, \dots, z_{0n}]$$

where $k^* = [k(X^*, X_1), \dots, k(X^*, X_n)]^T$

$$D^{-1} U^T (U D)^T U^T$$

$$= D U^T$$

$$\Rightarrow (D U^T)^T$$

$$\mathcal{D}'_{train} = \{(z_1, y_1), \dots, (z_n, y_n)\} \cup (1)$$

Then we compute:

1. Form \mathbb{K} on \mathcal{D}_{train}
2. Get UD
3. Pass $Z_q = (UD)[1 : q]$ to train \hat{f}
4. Form Z_q^* for all X_1^*, \dots, X_{n^*}
5. Get $\hat{Y}_{test} = \hat{f}(Z_q^*) \in \mathbb{R}^{n^*}$

KPCA: SUMMARY

$$(\mathcal{F}, \mathcal{F}^*)$$

$$f(x), \int f, \nabla f|_x$$

$$\forall \phi \in \mathcal{F}^* \Rightarrow \forall f \in \mathcal{F} \quad \phi f \in \mathbb{R} \quad \forall f \in \mathcal{F}$$

$$\sup_{\|f\|=1} \|\phi f\| \leq 1 \quad \exists \phi \in \mathcal{F}^* \text{ s.t. } \phi f = f(x)$$

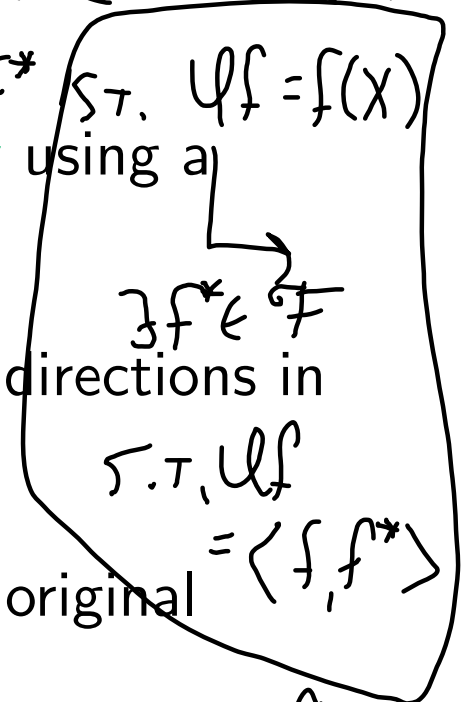
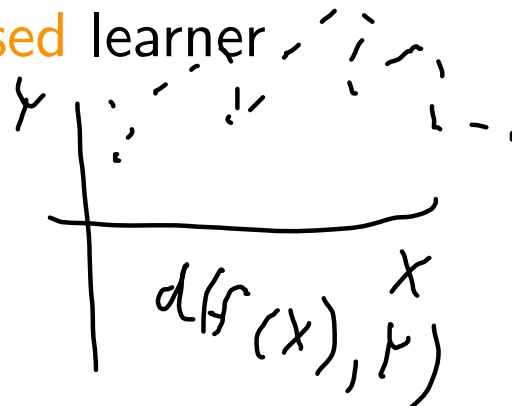
Kernel PCA seeks to generalize the notion of **similarity** using a kernel map

This can be interpreted as finding **smooth**, orthogonal directions in a RKHS

"SPLINE MODELS FOR OBSERVATIONAL DATA"
WAKABA

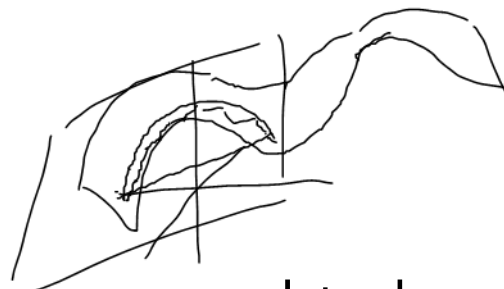
This can allow us to start picking up nonlinear (in the original feature space) aspects of our data

This new **representation** can be passed to a supervised method to form a **semisupervised** learner



\mathcal{F} RKHS

LAPLACIAN EIGENMAPS



In order to use the intuitive distance, we need to know the **geometry** of the data. This needs to be estimated.

We can get an estimate of the distance in the unknown geometry that the data come from (known as a **manifold**) by altering the usual Euclidean distance.

Some notes:

- The name **Laplacian Eigenmaps** comes from getting the **eigenvector** decomposition of the **Laplacian** restricted to the manifold (which is the second derivative version of the gradient).
- If the manifold is smooth, then **local** Euclidean distance is an approximation to the distance on the manifold.

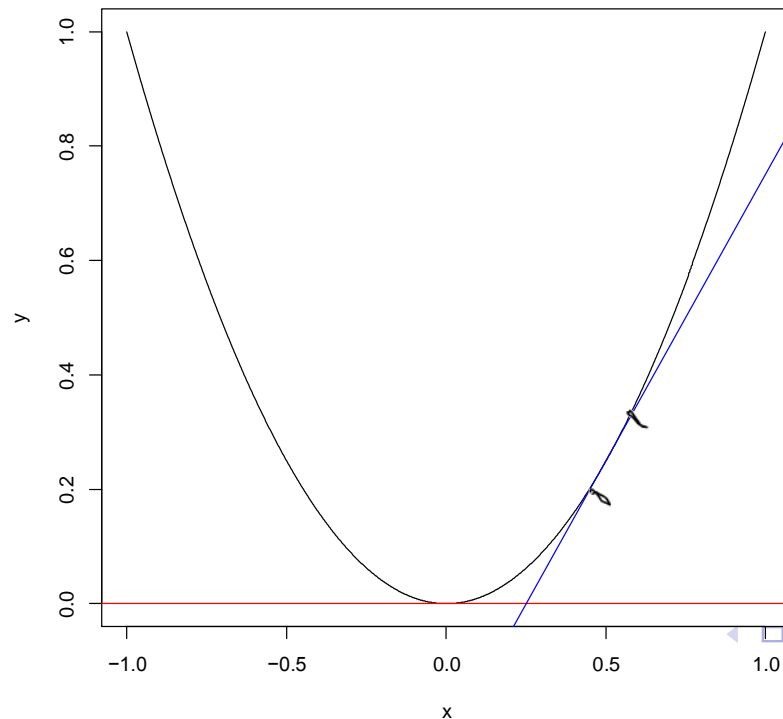
WHAT IS A MANIFOLD?

How good of an approximation is Euclidean distance?

This question is equivalent to how asking: how quickly does the **tangent** (space) change?

In 1-D, the tangent space is just the first derivative at that point:

$$f(x) = x^2 \Rightarrow f'(x) = 2x.$$



WHAT IS A MANIFOLD?

Therefore, the quality of the (local) Euclidean distance, depends on the **second derivative**

(ie: how fast does the first derivative change?)

In higher dimensions, the second derivative is known as the **Laplacian**:

$$\sum_j \frac{\partial^2 f}{\partial x_j^2}$$

(Note: This is also known as the **divergence** of the gradient)

WHAT ARE LAPLACIAN EIGENMAPS, THEN?

Imagine the operator \mathbb{L} that performs this operation:

$$Ax = b \qquad \mathbb{L}f = \sum_j \frac{\partial^2 f}{\partial x_j^2}$$

Then \mathbb{L} is the **Laplacian**, mapping a function to the divergence of its gradient

Key Idea: We can get the eigenvectors/eigenvalues of \mathbb{L} . Analogously to PCA, we can now do inference with these eigenvectors.


NOTE: There is a substantial overlap with KPCA, the difference being the centering of K and the **row sum** versus **column sum** normalization

LAPLACIAN EIGENMAPS

Collect data: X_1, \dots, X_n where $X_i \in \mathbb{R}^p$.

1. Form the distance matrix $\Delta_{ij} = \|X_i - X_j\|_2^2$.

2. Compute

$$M = (V, E)$$


$$K = \exp\left(-\frac{\Delta}{\gamma}\right)$$

3. Form the Laplacian $L = I - M^{-1}K$,

EIGENVECTOR OF ALL ONES

$$M = \text{diag}(\text{rowSums}(K))$$

↳ DEGREE

4. Compute the spectrum: $L = U\Sigma U^T$.

5. Return U_d , where U_d corresponds to the smallest d (nontrivial) eigenvalues of L

(Note that the eigenvectors of L and $M^{-1}K$ are the same but the order of the eigenvalues are reversed)

DEEPER INVESTIGATION

1. Form the distance matrix Δ .

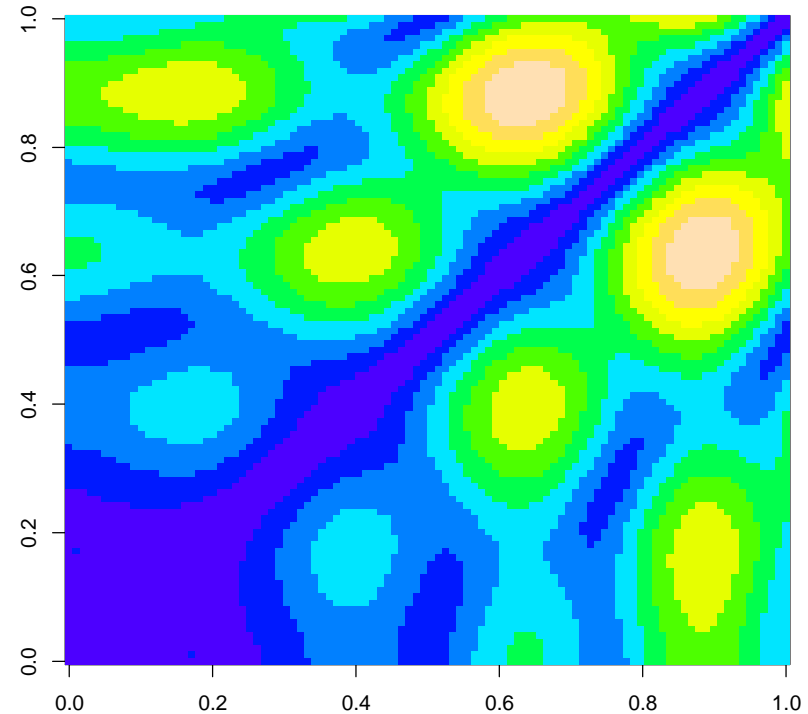
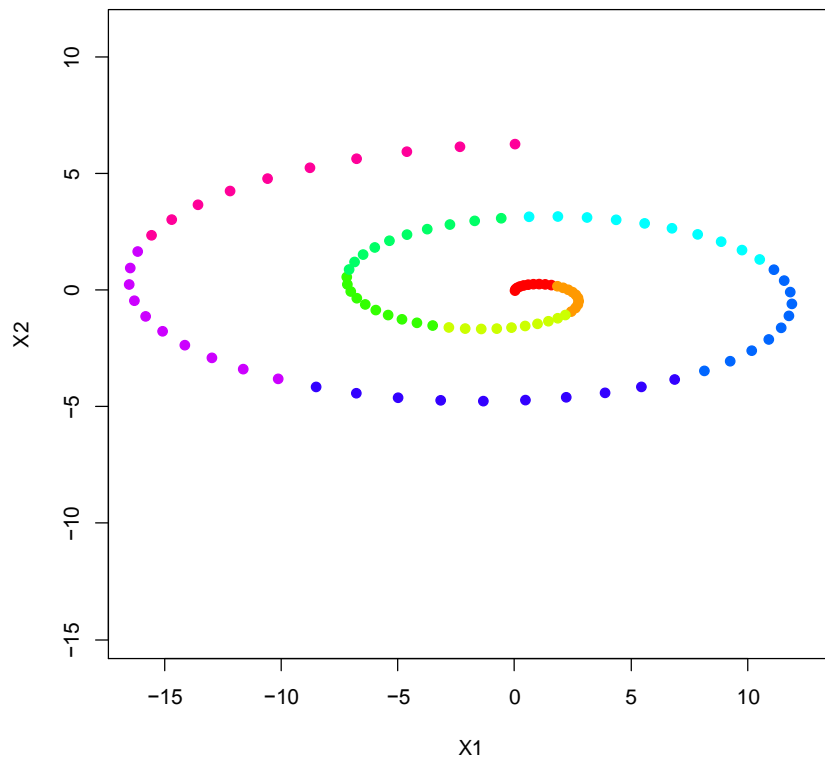


FIGURE: If we think about the center as 0 and the last blue circle as 1, then each entry the plot on the Right is the Euclidean distance between each data point on the plot on the Left (that is, Δ). The color on the Right plot goes from purple (small distance) to beige/pink (large distance).

DEEPER INVESTIGATION

```
Delta = as.matrix(dist(X,diag=TRUE,upper=TRUE))  
  
image(Delta,col=topo.colors(10))
```

DEEPER INVESTIGATION

2. Exponentiate $-\Delta/\gamma$ to form \mathbb{K} for some fixed γ .

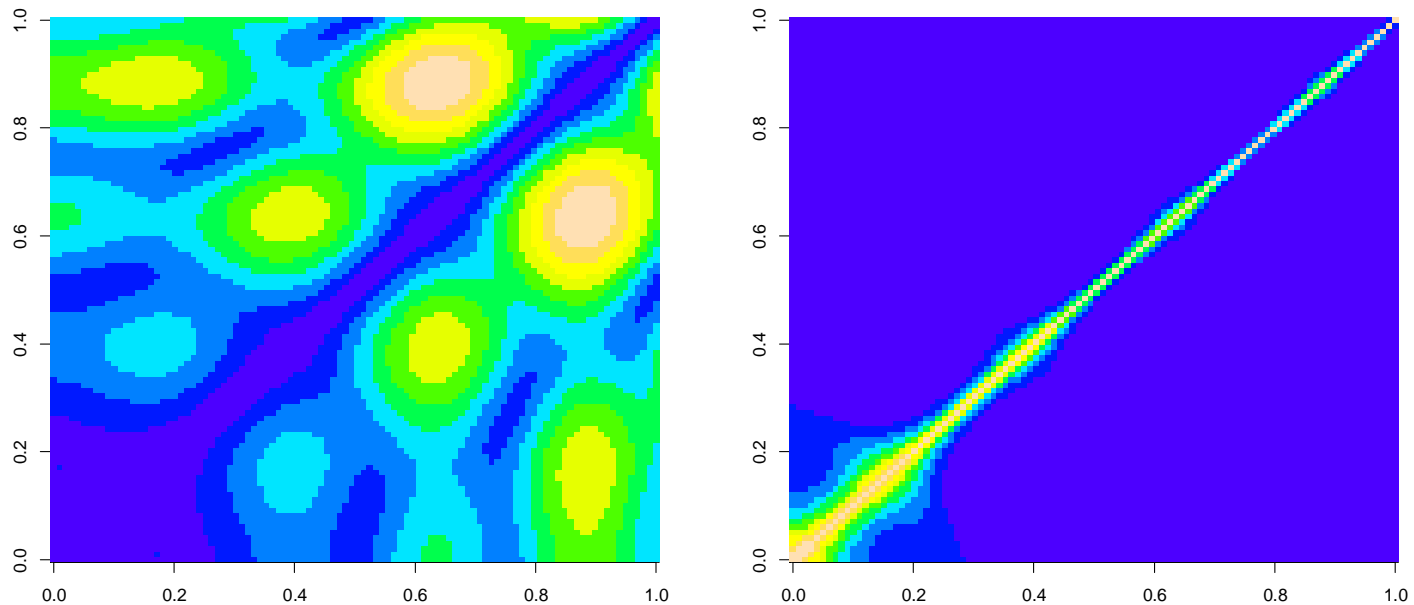
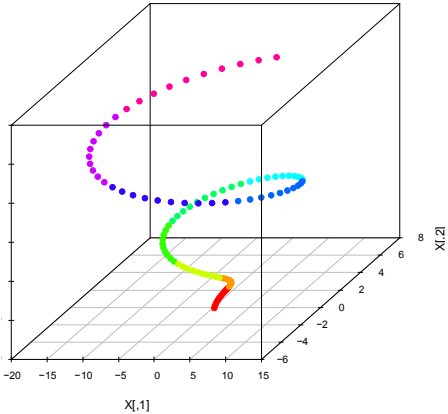


FIGURE: The Left plot is Δ and the Right plot is \mathbb{K} for $\gamma = 0.95$.

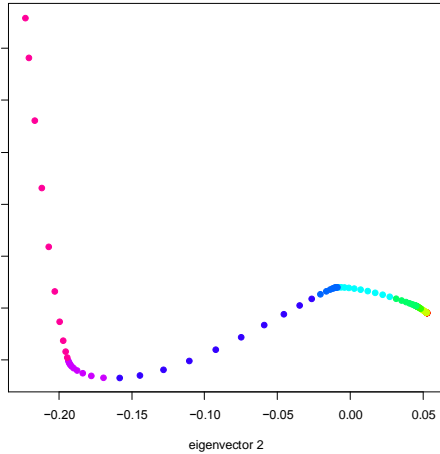
DEEPER INVESTIGATION

```
gamma = 0.95  
Wgamma = exp(-Delta/gamma)  
image(Wgamma,col=topo.colors(10))
```

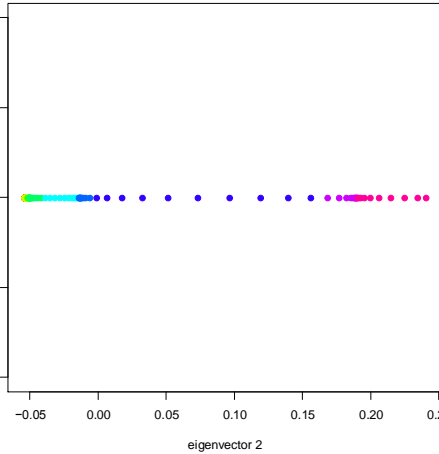

SPIRAL IN \mathbb{R}^3



Original data

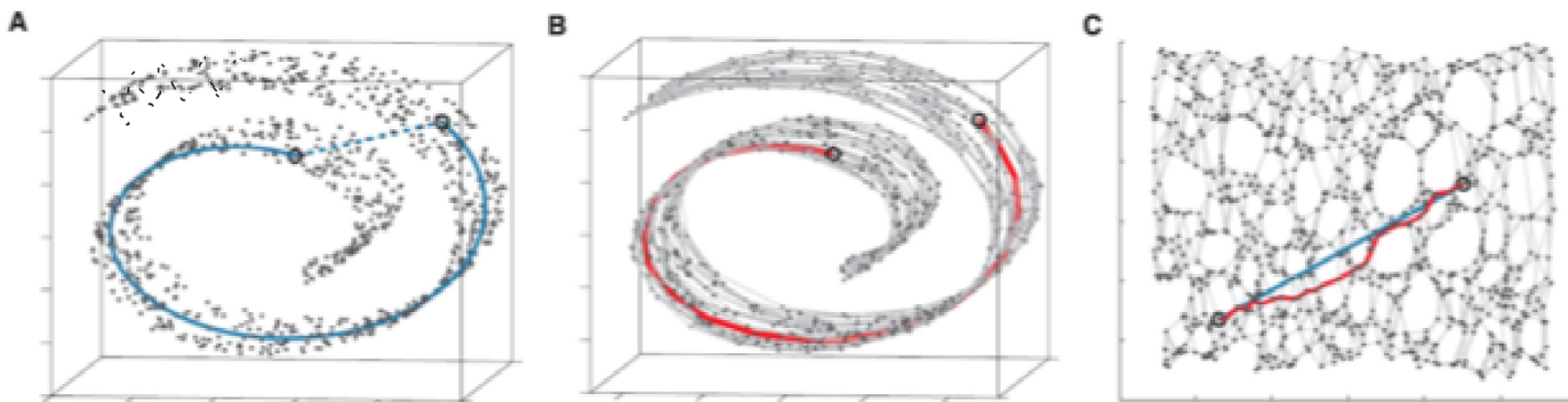


1st & 2nd nontrivial eigenvectors



1-dimensional

LOCAL EUCLIDEAN DISTANCE APPROXIMATES THE GEODESIC



The red line is the local Euclidean path between the two points, while the blue line is the path along the manifold.

James, Witten, Hastie, Tibshirani (2013)