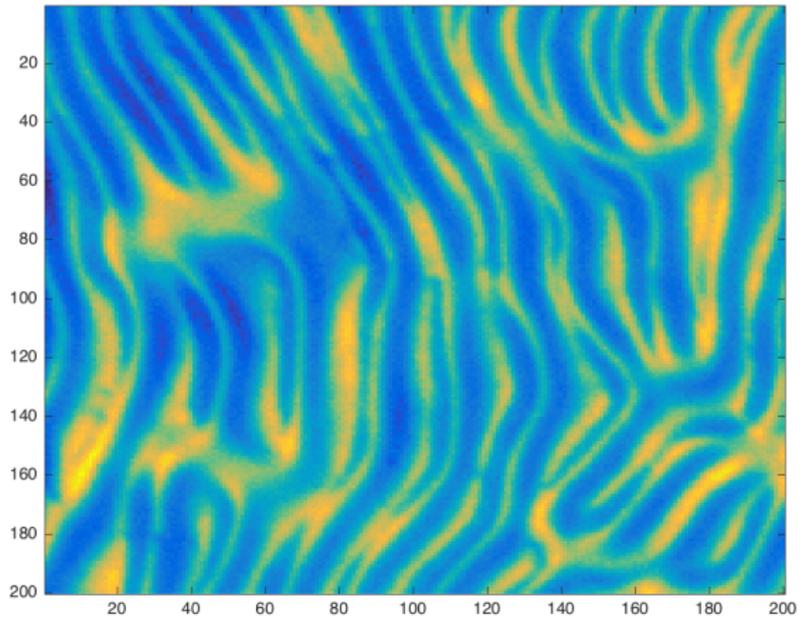


The Mathematics of Manifold Learning

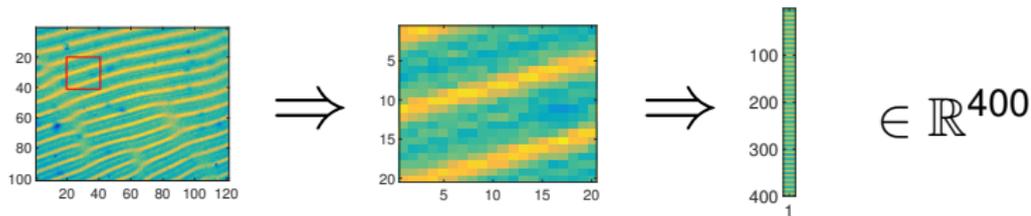
Tyrus Berry
George Mason University

April 6, 2019

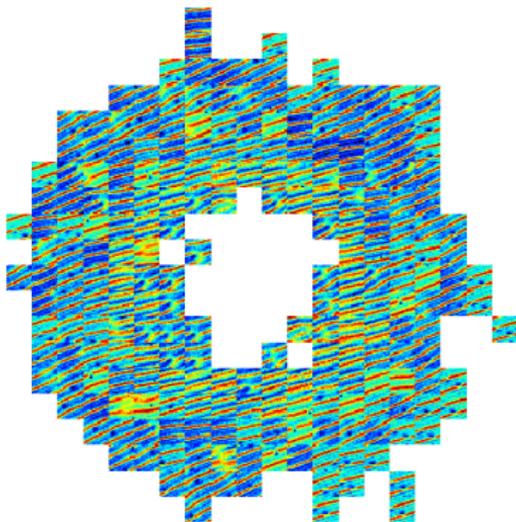
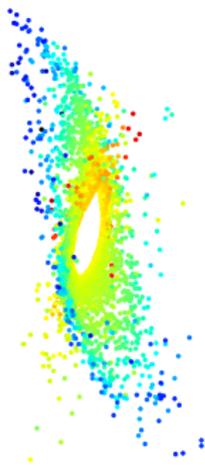
MOTIVATING EXAMPLE: NEMATIC LIQUID CRYSTAL



FINDING HIDDEN STRUCTURE IN DATA



The sub-image geometry:



OUTLINE

Lessons:

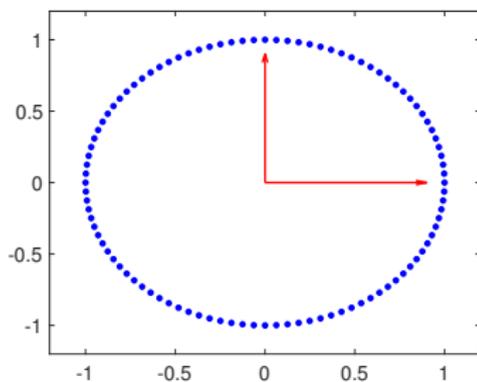
- ▶ **Dimensionality:** Intrinsic vs. Extrinsic
- ▶ **Nonlinearity:** Fourier Basis
- ▶ **Non-uniformity:** Respect the density

Challenges:

- ▶ Curse-of-dimensionality (intrinsic)
- ▶ Extrapolation

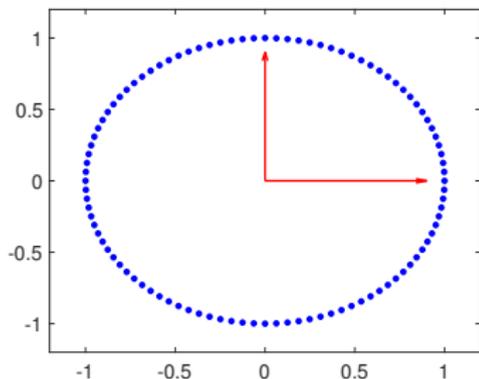
INTRINSIC VS. EXTRINSIC DIMENSION

100 points on a Circle



θ	x	y
0.0628	0.9980	0.0628
0.1257	0.9921	0.1253
0.1885	0.9823	0.1874
0.2513	0.9686	0.2487
0.3142	0.9511	0.3090
0.3770	0.9298	0.3681
0.4398	0.9048	0.4258
0.5027	0.8763	0.4818
⋮	⋮	
6.0319	0.9686	-0.2487
6.0947	0.9823	-0.1874
6.1575	0.9921	-0.1253
6.2204	0.9980	-0.0628
6.2832	1.0000	-0.0000

INTRINSIC VS. EXTRINSIC DIMENSION



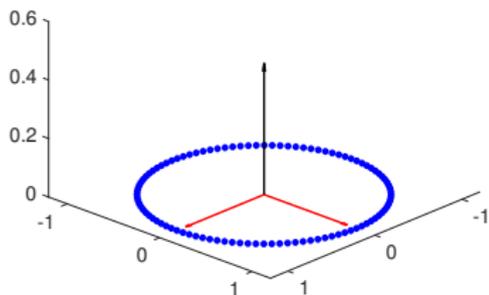
- ▶ Intrinsic Dimension = 1

$$\theta_i = 2\pi \frac{i}{100}$$

- ▶ Extrinsic Dimension = 2

$$(x_i, y_i) = (\cos(\theta_i), \sin(\theta_i))$$

INTRINSIC VS. EXTRINSIC DIMENSION



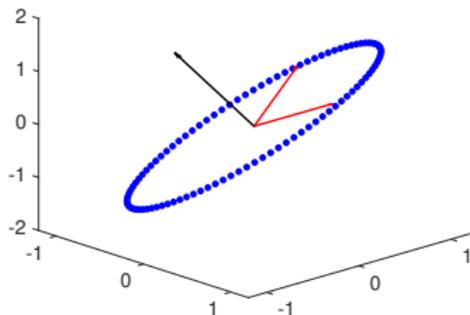
- ▶ Intrinsic Dimension = 1

$$\theta_i = 2\pi \frac{i}{100}$$

- ▶ Extrinsic Dimension = 3

$$(x_i, y_i, z_i) = (\cos(\theta_i), \sin(\theta_i), 0)$$

INTRINSIC VS. EXTRINSIC DIMENSION



- ▶ Intrinsic Dimension = 1

$$\theta_i = 2\pi \frac{i}{100}$$

- ▶ Extrinsic Dimension = 3

$$x_i = \cos(\theta_i)$$

$$y_i = \sin(\theta_i)$$

$$z_i = x_i + y_i$$

INTRINSIC VS. EXTRINSIC DIMENSION

- ▶ Intrinsic Dimension = 1

$$\theta_i = 2\pi \frac{i}{100}$$

- ▶ Extrinsic Dimension = $2 + n$

$$\begin{aligned} x_i &= \cos(\theta_i) \\ y_i &= \sin(\theta_i) \\ z_i^1 &= a_1 x_i + b_1 y_i \\ &\vdots \\ z_i^n &= a_n x_i + b_n y_i \end{aligned} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ a_1 & a_2 \\ \vdots & \vdots \\ a_n & b_n \end{bmatrix} \begin{bmatrix} \cos(\theta_i) \\ \sin(\theta_i) \end{bmatrix} = A \begin{bmatrix} \cos(\theta_i) \\ \sin(\theta_i) \end{bmatrix}$$

A is a $(n + 2) \times 2$ matrix

SOLUTION: LINEAR ALGEBRA!

- ▶ Hidden Data, $[\theta_1, \theta_2, \theta_3, \dots, \theta_N]$
- ▶ Ideal Representation, $x_i = \cos(\theta_i), y_i = \sin(\theta_i)$

$$X = \begin{bmatrix} x_1 & x_2 & x_3 & \cdots & x_N \\ y_1 & y_2 & y_3 & \cdots & y_N \end{bmatrix}$$

- ▶ Given Data: $Y = AX$

$$Y = \begin{bmatrix} x_1 & x_2 & x_3 & \cdots & x_N \\ y_1 & y_2 & y_3 & \cdots & y_N \\ a_1 x_1 + b_1 y_1 & a_1 x_2 + b_1 y_2 & a_1 x_3 + b_1 y_3 & \cdots & a_1 x_N + b_1 y_N \\ a_2 x_1 + b_2 y_1 & a_2 x_2 + b_2 y_2 & a_2 x_3 + b_2 y_3 & \cdots & a_2 x_N + b_2 y_N \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_n x_1 + b_n y_1 & a_n x_2 + b_n y_2 & a_n x_3 + b_n y_3 & \cdots & a_n x_N + b_n y_N \end{bmatrix}$$

- ▶ Rows of Y are linearly dependent!

SOLUTION: LINEAR ALGEBRA!

- ▶ Given data $Y = AX$ where both A and X are unknown
- ▶ Linear dependence means the rows, Y_i , are redundant:

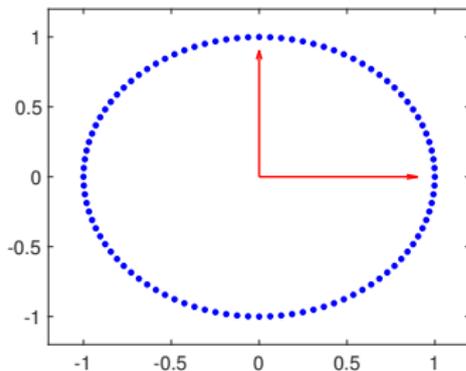
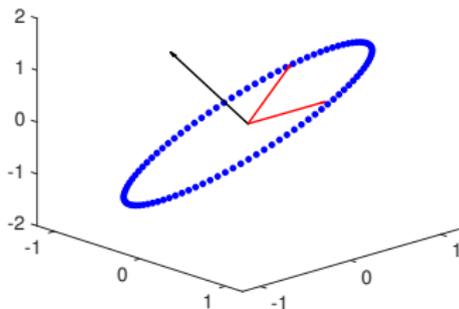
$$\vec{c}^\top Y = c_1 Y_1 + c_2 Y_2 + \dots + c_n Y_n = \vec{0}$$

- ▶ There exists $\vec{c} = (c_1, \dots, c_n) \neq 0$ such that $\vec{c}^\top Y = \vec{0}$
- ▶ $\vec{c}^\top Y = \vec{0}$ if and only if $\vec{c}^\top YY^\top \vec{c} = \vec{0}^\top \vec{0} = 0$
- ▶ So \vec{c} is eigenvector of YY^\top with eigenvalue zero

PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ Compute the eigenvectors/values of $YY^T = U\Lambda U^T$
- ▶ Sort the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$
- ▶ Eigenvalue ≈ 0 represent linear redundancies
- ▶ **Principal Components:** Eigenvectors u_i with largest λ_i
- ▶ Choose $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_p$ corresponding to $\lambda_1, \dots, \lambda_p$
- ▶ Form the projection matrix $P = [\vec{u}_1 \ \vec{u}_2 \ \dots \ \vec{u}_p]$
- ▶ Remove redundancies: $\tilde{X} = PY$

PRINCIPAL COMPONENT ANALYSIS (PCA)

 Y \Rightarrow $\tilde{X} = PY$ 

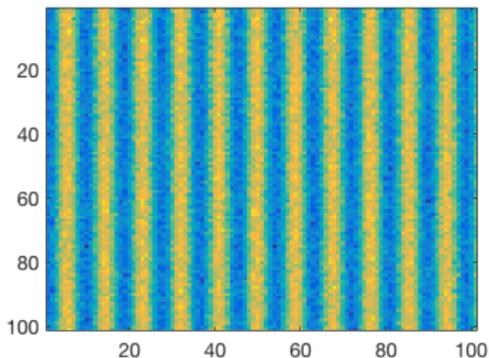
PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ Matrix times *intrinsic* data \Rightarrow *extrinsic* redundancy
- ▶ These *linear* redundancies are easy to remove
- ▶ PCA projects the data to remove redundancy
- ▶ Does this really happen?

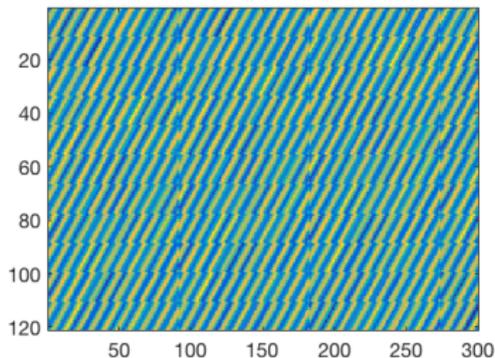
DOES THIS REALLY HAPPEN?

Consider 11×11 subimages from a pattern:

Original Image

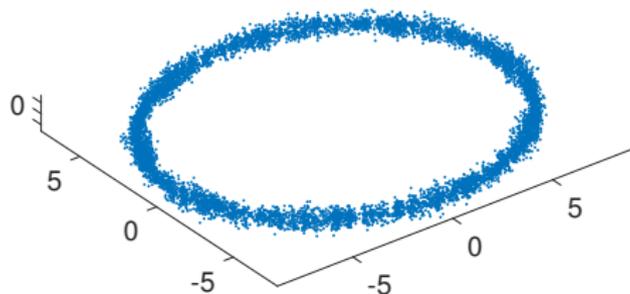


Vectorized Subimages

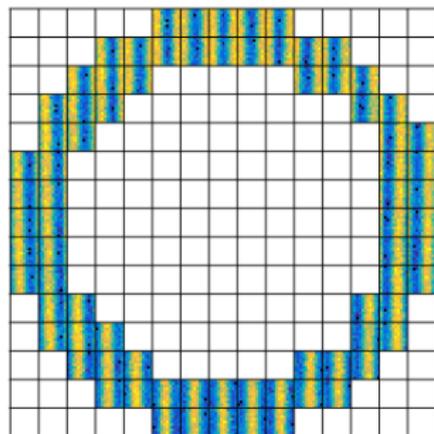


DOES THIS REALLY HAPPEN?

PCA Coordinates

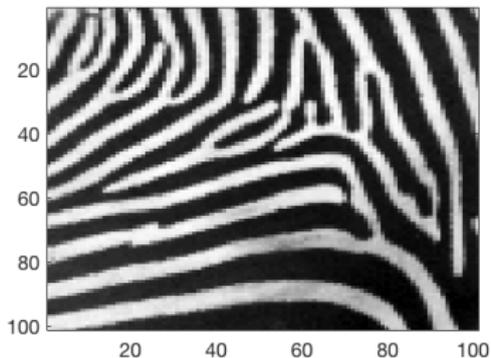


Subimage Coordinates

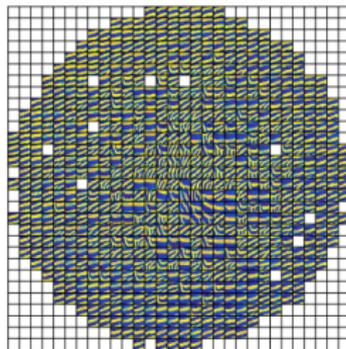
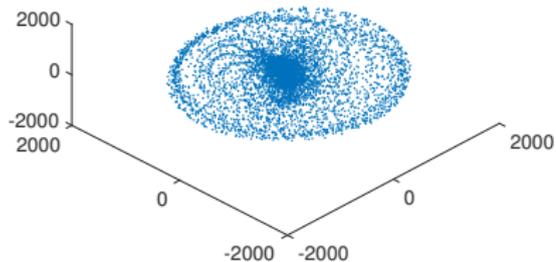


DOES THIS REALLY HAPPEN?

Zebra Stripes

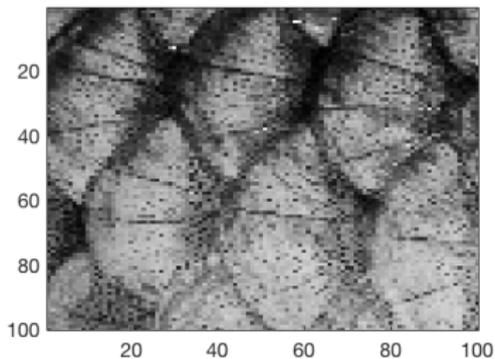


PCA Coordinates

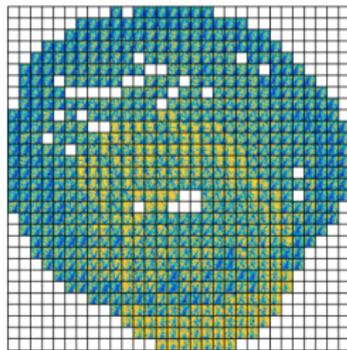
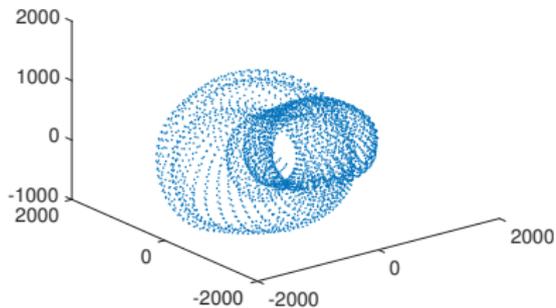


DOES THIS REALLY HAPPEN?

Fish Scales

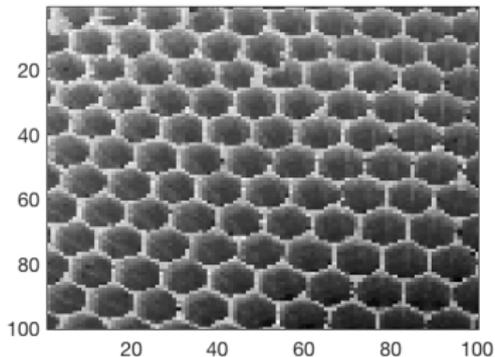


PCA Coordinates

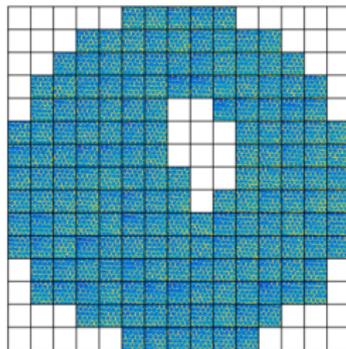
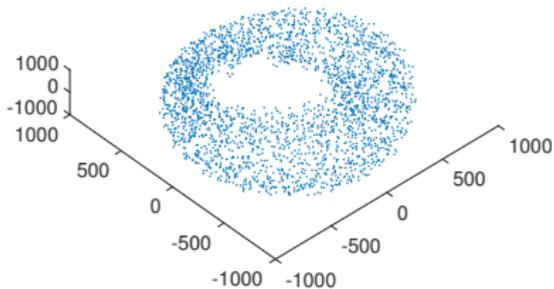


DOES THIS REALLY HAPPEN?

Honeycomb



PCA Coordinates



PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ *Linear* redundancies are easy to remove

$$c_1 Y_1 + c_2 Y_2 + \cdots c_n Y_n = \vec{0}$$

PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ *Linear* redundancies are easy to remove

$$c_1 Y_1 + c_2 Y_2 + \cdots c_n Y_n = \vec{0}$$

- ▶ PCA projects the data to remove redundancy

PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ *Linear* redundancies are easy to remove

$$c_1 Y_1 + c_2 Y_2 + \dots + c_n Y_n = \vec{0}$$

- ▶ PCA projects the data to remove redundancy
- ▶ What about **nonlinear** redundancies?

$$F(Y_1, Y_2, \dots, Y_n) = \vec{0}$$

PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ *Linear* redundancies are easy to remove

$$c_1 Y_1 + c_2 Y_2 + \cdots c_n Y_n = \vec{0}$$

- ▶ PCA projects the data to remove redundancy
- ▶ What about **nonlinear** redundancies?

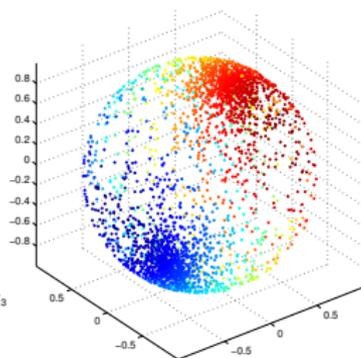
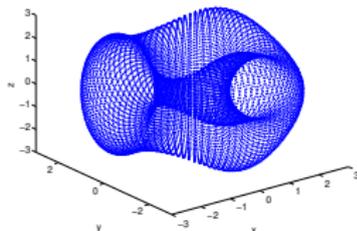
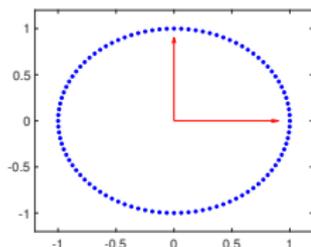
$$F(Y_1, Y_2, \dots, Y_n) = \vec{0}$$

- ▶ Example, Circle: $Y_1 = \cos(\theta)$, $Y_2 = \sin(\theta)$

$$F(Y_1, Y_2) = Y_1^2 + Y_2^2 - 1 = \vec{0}$$

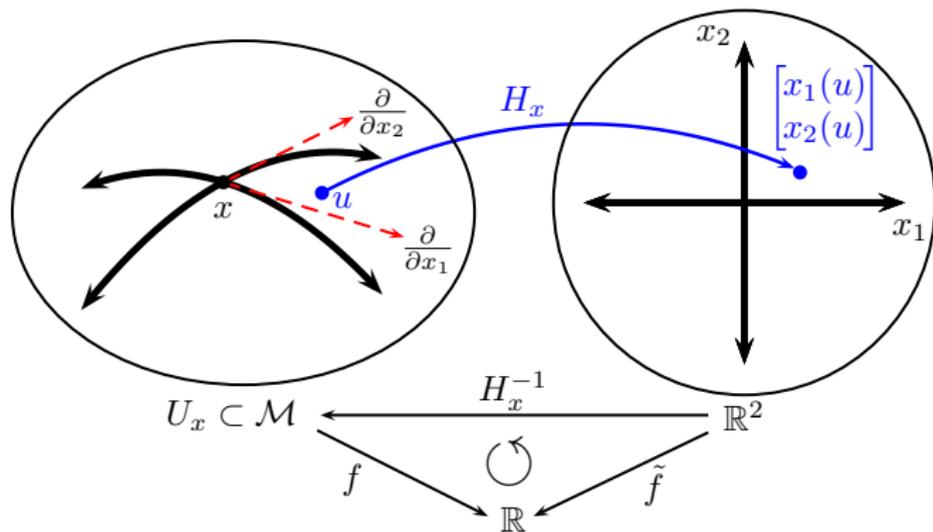
MANIFOLD LEARNING

A *manifold* \mathcal{M} is a topological space that is *locally* Euclidean.



MANIFOLD LEARNING

Around each point $x \in \mathcal{M}$ we have an open neighborhood $U_x \subset \mathcal{M}$ and a homeomorphism $H_x : U_x \rightarrow \mathbb{R}^m$



MANIFOLD LEARNING

- ▶ When does a nonlinear redundancy define a manifold?

$$\mathcal{M} = \{\vec{y} \mid F(y_1, y_2, \dots, y_n) = \vec{a}\} \subset \mathbb{R}^n$$

MANIFOLD LEARNING

- ▶ When does a nonlinear redundancy define a manifold?

$$\mathcal{M} = \{\vec{y} \mid F(y_1, y_2, \dots, y_n) = \vec{a}\} \subset \mathbb{R}^n$$

- ▶ Need to be able to solve for the last $n - m$ variables:

$$\begin{aligned}\vec{a} &= F(y_1, \dots, y_m, y_{m+1}, \dots, y_n) \\ &= F(y_1, \dots, y_m, G_1(y_1, \dots, y_m), G_2(y_1, \dots, y_m), \dots, G_{n-m}(y_1, \dots, y_m))\end{aligned}$$

MANIFOLD LEARNING

- ▶ When does a nonlinear redundancy define a manifold?

$$\mathcal{M} = \{\vec{y} \mid F(y_1, y_2, \dots, y_n) = \vec{a}\} \subset \mathbb{R}^n$$

- ▶ Need to be able to solve for the last $n - m$ variables:

$$\begin{aligned}\vec{a} &= F(y_1, \dots, y_m, y_{m+1}, \dots, y_n) \\ &= F(y_1, \dots, y_m, G_1(y_1, \dots, y_m), G_2(y_1, \dots, y_m), \dots, G_{n-m}(y_1, \dots, y_m))\end{aligned}$$

- ▶ **Implicit Function Theorem:** If the Jacobian matrix $DF(\vec{y})$ is full rank then the functions G_1, \dots, G_{n-m} exist near \vec{y}

MANIFOLD LEARNING

- ▶ When does a nonlinear redundancy define a manifold?

$$\mathcal{M} = \{\vec{y} \mid F(y_1, y_2, \dots, y_n) = \vec{a}\} \subset \mathbb{R}^n$$

- ▶ Need to be able to solve for the last $n - m$ variables:

$$\begin{aligned}\vec{a} &= F(y_1, \dots, y_m, y_{m+1}, \dots, y_n) \\ &= F(y_1, \dots, y_m, G_1(y_1, \dots, y_m), G_2(y_1, \dots, y_m), \dots, G_{n-m}(y_1, \dots, y_m))\end{aligned}$$

- ▶ **Implicit Function Theorem:** If the Jacobian matrix $DF(\vec{y})$ is full rank then the functions G_1, \dots, G_{n-m} exist near \vec{y}
- ▶ **Sard's Theorem:** If F is smooth, then for *almost every* \vec{a} , the Jacobian $DF(\vec{y})$ is full rank for all $\vec{y} \in \mathcal{M}$

MANIFOLD LEARNING

- ▶ When does a nonlinear redundancy define a manifold?

$$\mathcal{M} = \{\vec{y} \mid F(y_1, y_2, \dots, y_n) = \vec{a}\} \subset \mathbb{R}^n$$

MANIFOLD LEARNING

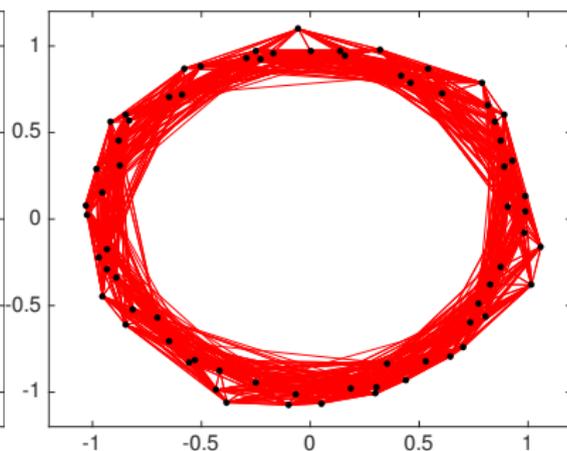
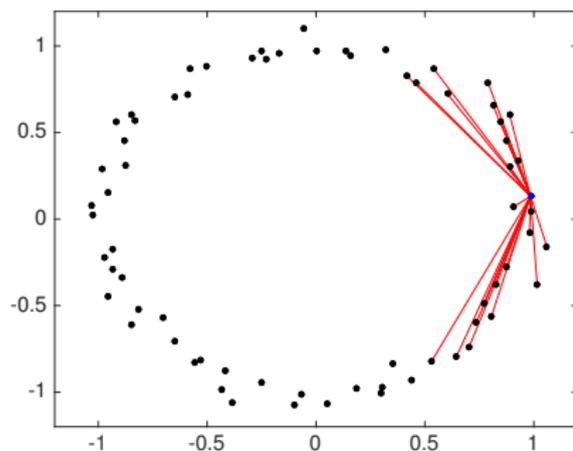
- ▶ When does a nonlinear redundancy define a manifold?

$$\mathcal{M} = \{\vec{y} \mid F(y_1, y_2, \dots, y_n) = \vec{a}\} \subset \mathbb{R}^n$$

- ▶ When F is smooth, \mathcal{M} is a manifold for almost every \vec{a}

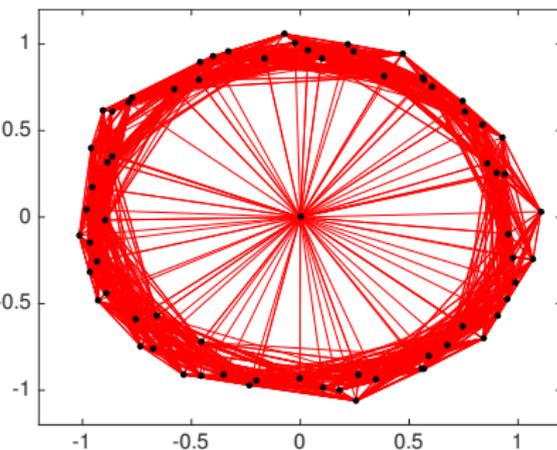
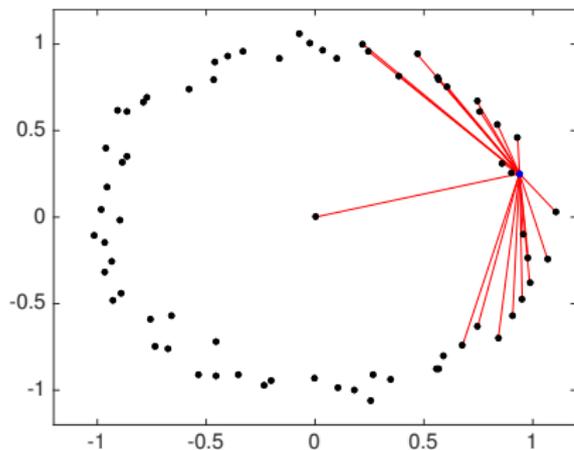
MANIFOLD \Rightarrow GRAPH

- ▶ Represent the **nonlinear** structure with a graph
- ▶ Locally Euclidean \Rightarrow Connect nearby points



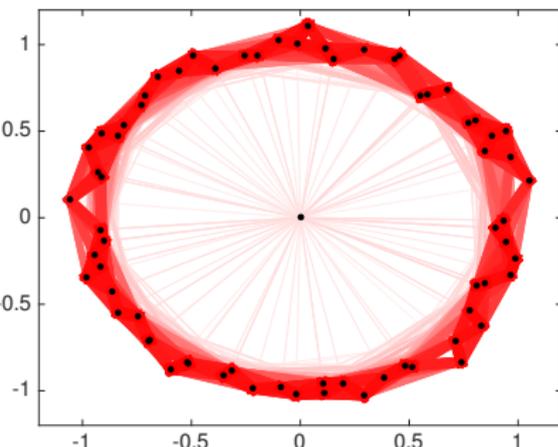
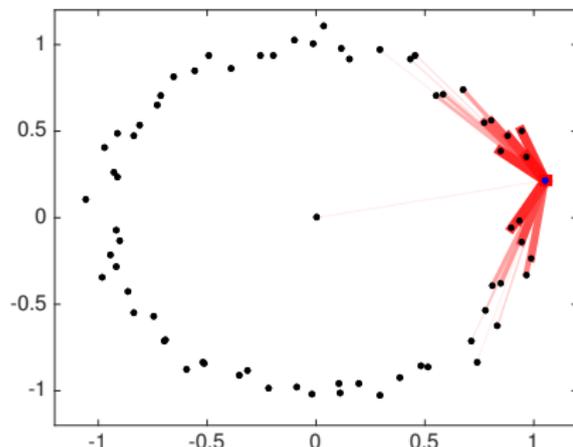
MANIFOLD \Rightarrow GRAPH

- **Problem:** Noise and outliers can lead to **bridging**



MANIFOLD \Rightarrow GRAPH

- ▶ To prevent bridging, edges weighted: $K_\delta(x, y) = e^{-\frac{\|x-y\|^2}{4\delta^2}}$
- ▶ **Theorem:** Graph encodes all nonlinear information



WHAT IS MANIFOLD LEARNING?

- ▶ **Manifold learning** \Leftrightarrow **Estimating Laplace Operator**
- ▶ Euclidean space:
 - ▶ Eigenfunctions of Δ are $e^{i\vec{\omega}\cdot\vec{x}}$
 - ▶ Basis for Fourier transform
- ▶ Unit circle:
 - ▶ Eigenfunctions of Δ are $e^{ik\theta}$
 - ▶ Basis for Fourier series
- ▶ **Theorem:** Eigenfunctions of Δ give the smoothest basis for square integrable functions on any manifold.

FINDING THE LAPLACIAN FROM DATA

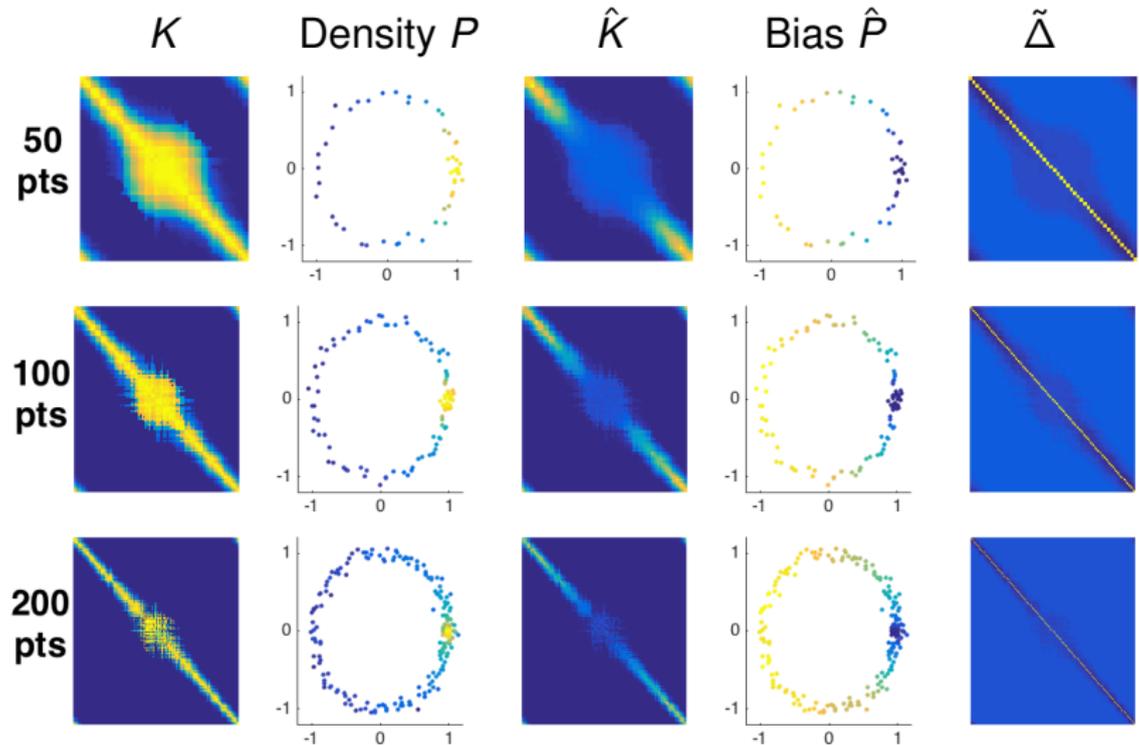
- ▶ We have converted our data set to a *weighted graph*
- ▶ Vertices \Rightarrow Data points $\{x_1, x_2, \dots, x_N\}$
- ▶ Edges \Rightarrow Pairs of nearest neighbors
- ▶ Edge Weights $\Rightarrow K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{4\epsilon}}$
- ▶ Represented as matrix $K_{ij} = K(x_i, x_j)$

DIFFUSION MAPS: THE KEY RESULT

1. Start with the matrix $K_{ij} = e^{-\frac{\|x_i - x_j\|^2}{4\epsilon}}$
2. Find the row sums $P_i = \sum_{j=1}^N K_{ij}$
3. Normalize the matrix $\hat{K}_{ij} = \frac{K_{ij}}{P_i P_j}$
4. Find the row sums again $\hat{P}_i = \sum_{j=1}^N \hat{K}_{ij}$
5. Markov Normalization $\tilde{K}_{ij} = \frac{\hat{K}_{ij}}{\hat{P}_i}$
6. Form the Laplacian matrix $\tilde{\Delta} = \frac{I - \tilde{K}}{\epsilon}$

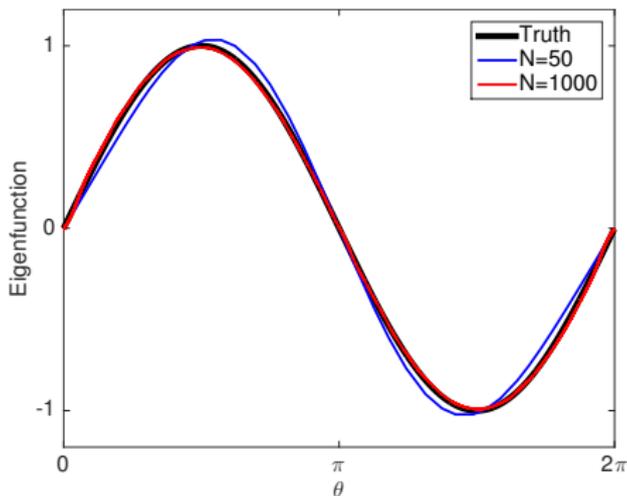
Theorem: As $N \rightarrow \infty$ and $\epsilon \rightarrow 0$ we have $\tilde{\Delta} \rightarrow \Delta$

DIFFUSION MAPS CONSTRUCTION



DIFFUSION MAPS CONSTRUCTION

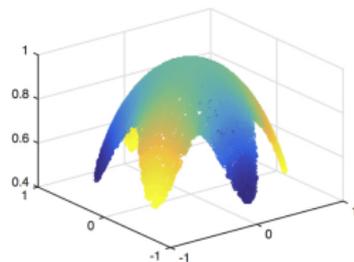
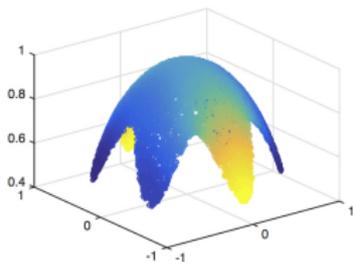
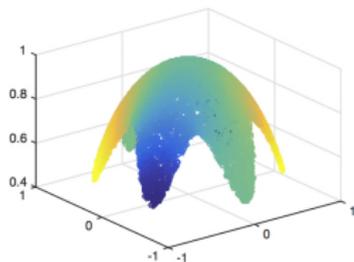
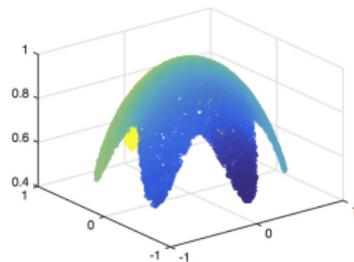
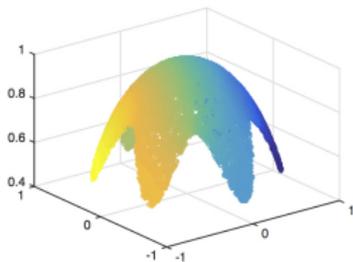
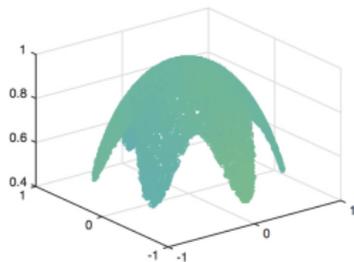
- ▶ $\tilde{\Delta}$ approximates the Laplacian Δ
- ▶ $\tilde{\Delta}$ encodes the geometry of the data
- ▶ Eigenvectors of $\tilde{\Delta}$ approximate eigenfunctions of Δ



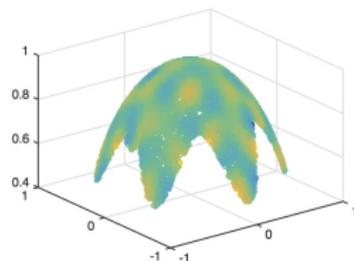
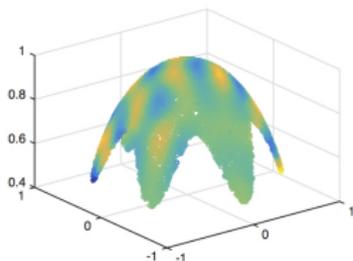
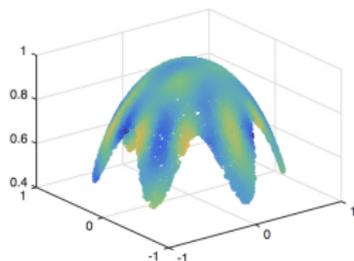
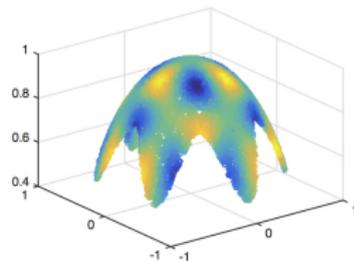
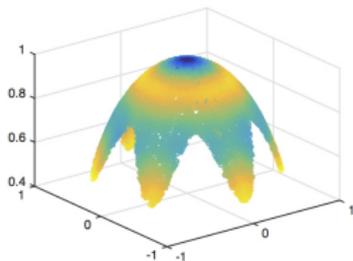
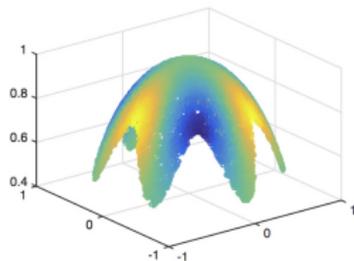
FOURIER BASIS ON MANIFOLDS

- ▶ Fourier functions $\sin(k\theta)$ are eigenfunctions of $\frac{d^2}{d\theta^2}$
- ▶ Eigenvectors of matrix $\tilde{\Delta}$ approximate eigenfunctions of Δ
- ▶ What is so great about these functions?
 - ▶ Smoothest possible functions on \mathcal{M}
 - ▶ $\varphi_0 = \text{constant}$
 - ▶ φ_1 contains a single oscillation
 - ▶ φ_j is smoothest function orthogonal to previous

FOURIER BASIS ON MANIFOLDS



FOURIER BASIS ON MANIFOLDS



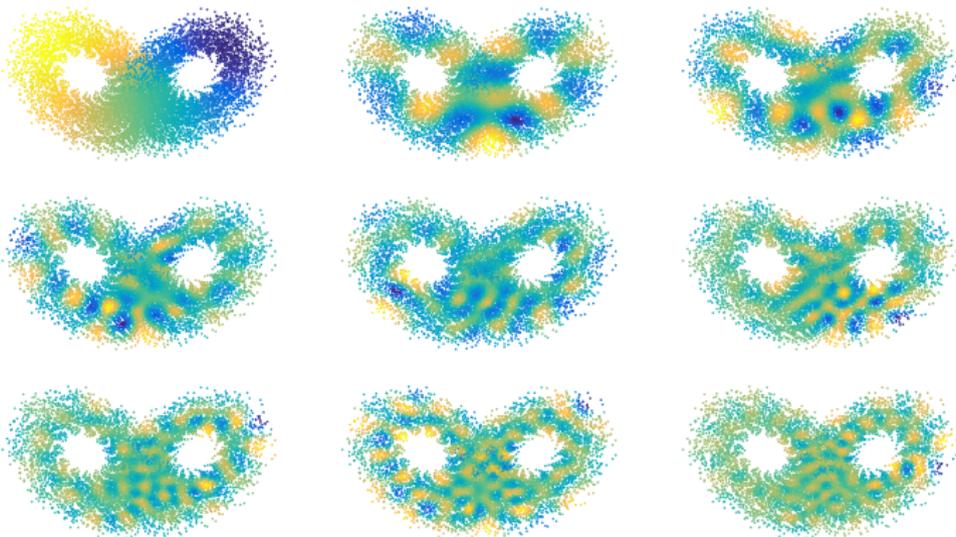
FORECASTING WITHOUT A MODEL

$$\begin{array}{ccc}
 \rho(x, t) & \xrightarrow{\text{Nonparametric Forecast}} & \rho(x, t + \tau) \\
 \downarrow \langle \rho, \varphi_j \rangle & & \uparrow \sum_j c_j \varphi_j \rho_{\text{eq}} \\
 \vec{c}(t) & \xrightarrow{A_{lj} \equiv \mathbb{E}[\langle \varphi_j, \mathcal{S} \varphi_l \rangle \rho_{\text{eq}}]} & \vec{c}(t + \tau) = A \vec{c}(t).
 \end{array}$$

- ▶ $\vec{c}(t)$ are the generalized Fourier coefficients of ρ
- ▶ Nonlinear dynamics become linear (matrix A) in this basis

MANIFOLD LEARNING \Rightarrow CUSTOM 'FOURIER' BASIS

- ▶ **Optimal basis:** Minimum variance $A_{lj} \equiv \mathbb{E}[\langle \varphi_l, S\varphi_l \rangle_q]$

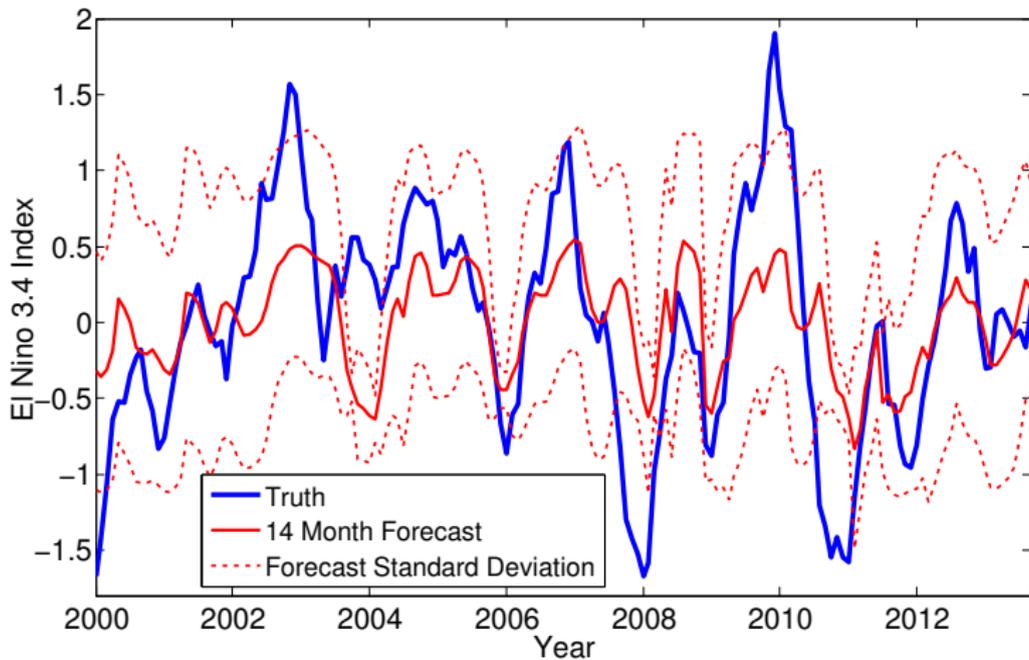


EXAMPLE: FORECASTING WITHOUT A MODEL

No Model

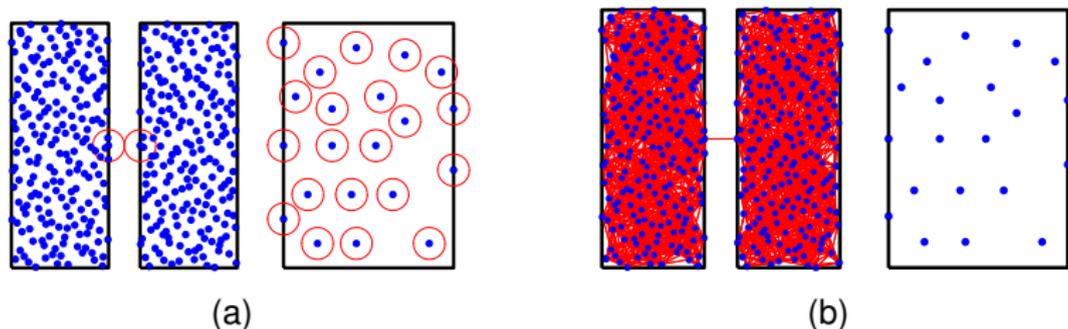
Perfect Model

EXAMPLE: FORECASTING EL NINO



NONUNIFORM DENSITY: FIXED BALLS

Black outlines indicate true clusters:

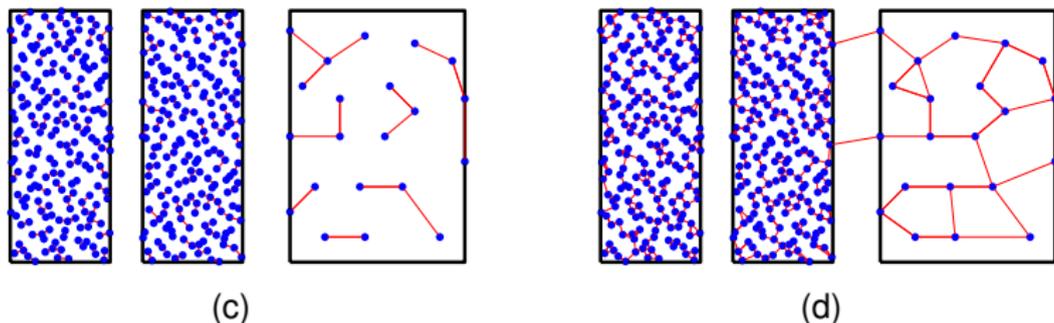


(a) Dense regions **bridged** before connecting sparse region

(b) Graph connecting all points with distance less than ϵ

$$\|x - y\| < \epsilon$$

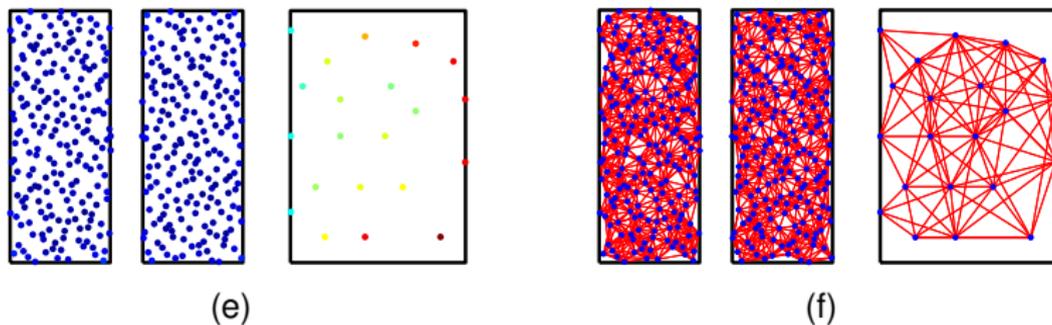
NONUNIFORM DENSITY: NEAREST NEIGHBORS (NN)



(c) Connect each point to its **nearest neighbor** (NN)

(d) Connect each point to its two nearest neighbors (2NN)

NONUNIFORM DENSITY: CKNN

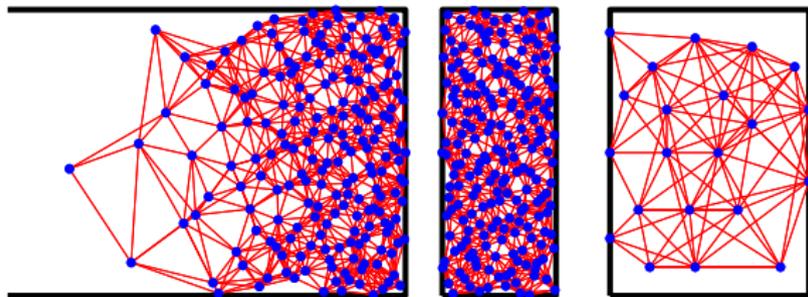


(e) Distance to 10-th nearest neighbor

(f) **Continuous k-Nearest Neighbors (CkNN)**

$$\frac{\|x - y\|}{\sqrt{\|x - \text{kNN}(x)\| \cdot \|y - \text{kNN}(y)\|}} < \delta$$

NONUNIFORM DENSITY: CONCLUSION



(h)

(h) Real data has sparse tails: More data = bigger gaps!

Theorem: NN fails even with infinite data. CkNN succeeds.

IMPROVED CLUSTERING USING CKNN

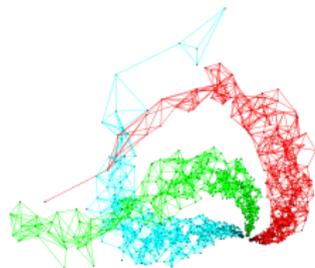
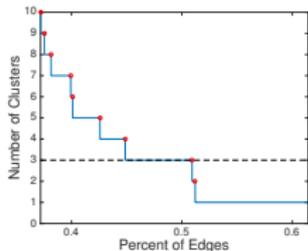
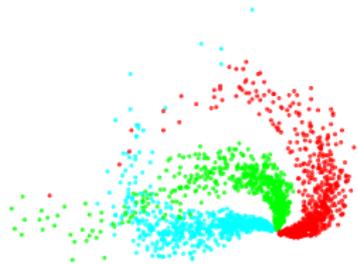
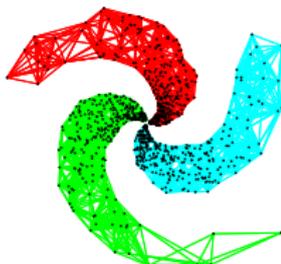
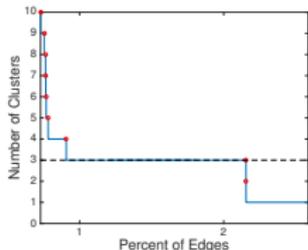
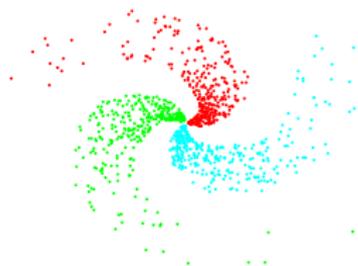
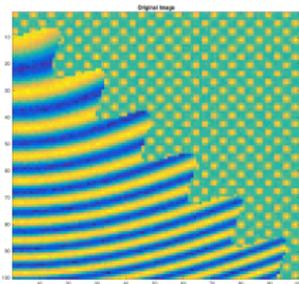


IMAGE SEGMENTATION

Original Image: Break into subimages



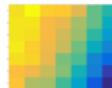
(a)



(b)



(c)



(d)



(e)



(f)



(g)



(h)



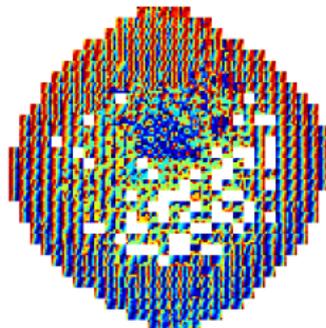
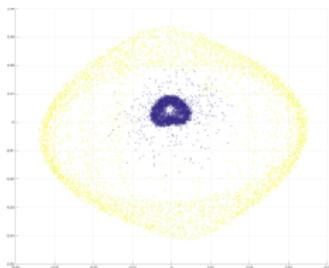
(i)

Images produced by Marilyn Vazquez.

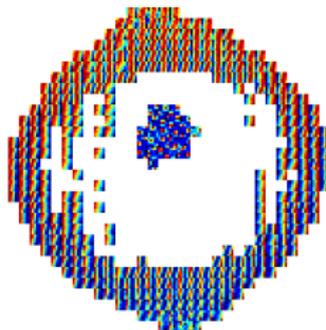
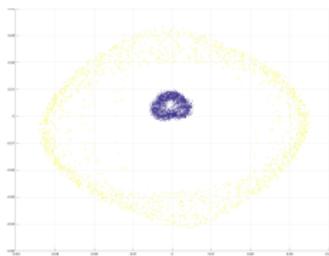
IMAGE SEGMENTATION

Clustering shown projected to two principal components

all points



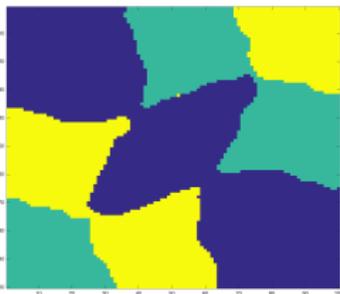
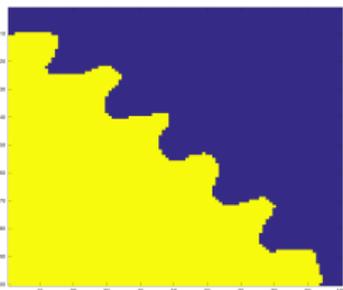
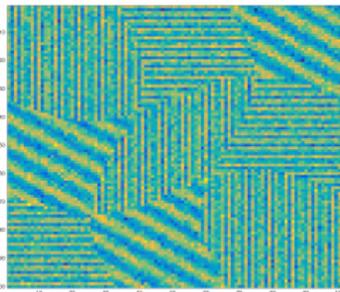
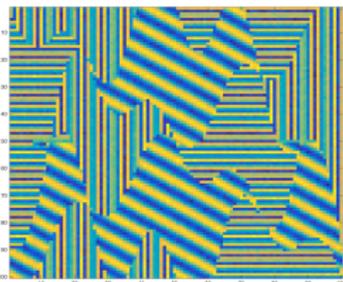
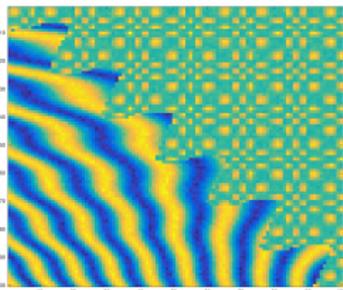
with low
density
points
removed



Images produced by Marilyn Vazquez.

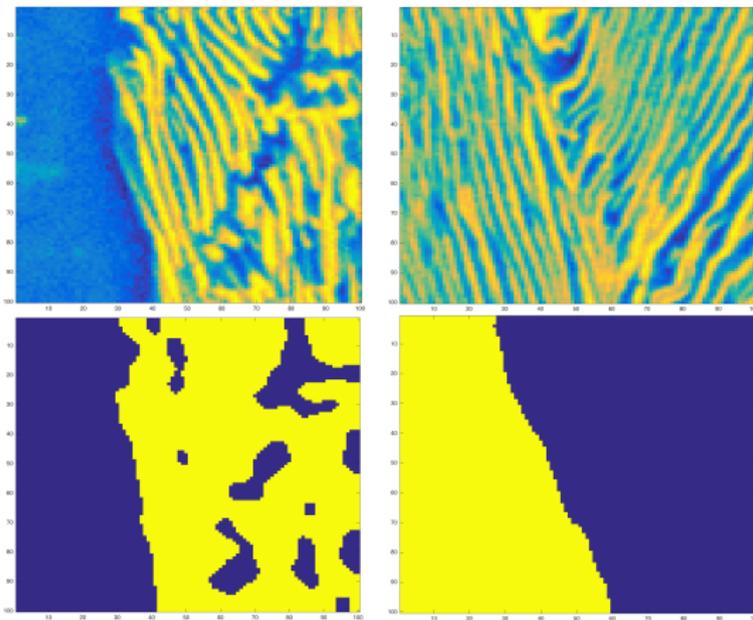
IMAGE SEGMENTATION

Results - synthetic images



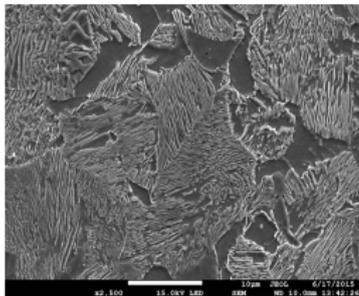
Images produced by Marilyn Vazquez.

IMAGE SEGMENTATION: REAL IMAGES

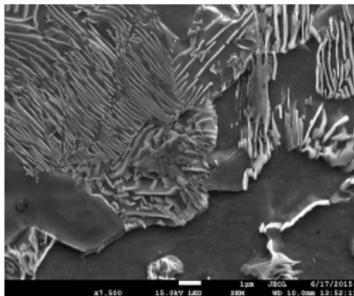


Images produced by Marilyn Vazquez.

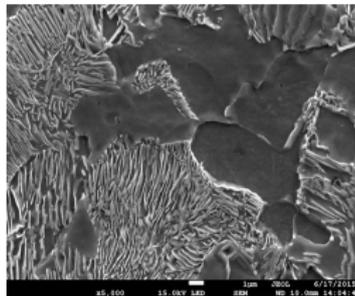
IMAGE SEGMENTATION: REAL IMAGES



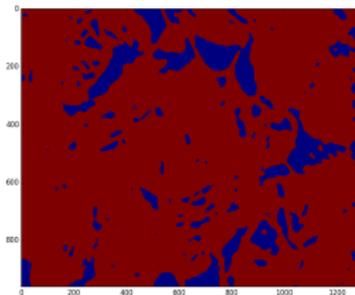
(g)



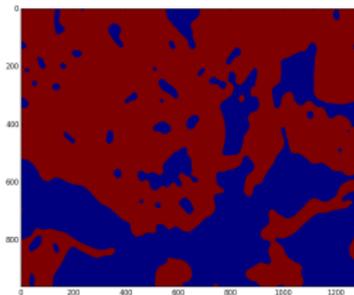
(h)



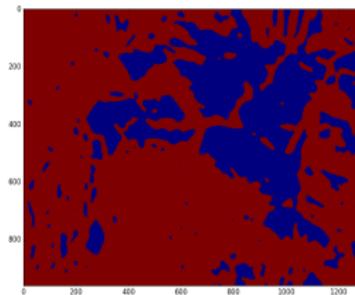
(i)



(j)



(k)



(l)

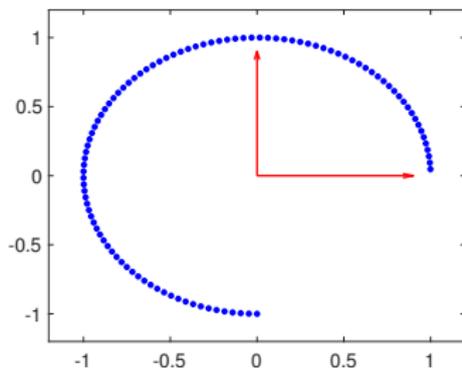
Original images by Mark R. Stoudt and Steve P. Mates. Analysis by Marilyn Vazquez.

CURSE-OF-(INTRINSIC)-DIMENSIONALITY

- ▶ **Try** to cut into independent components
- ▶ Otherwise math/stat says it is **impossible**
- ▶ Need **more/better** assumptions and/or questions
- ▶ **Better assumptions:** Smoothness
- ▶ **Better questions:** Feature of interest (supervised)

EXTRAPOLATION

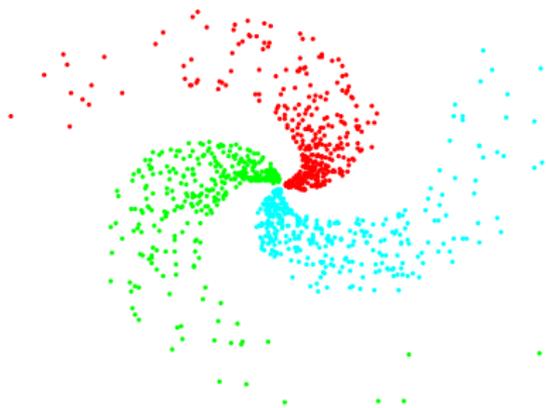
- ▶ Given only part of a structure recover the whole



- ▶ Need to exploit symmetry

EXTRAPOLATION

- ▶ Given only part of a structure recover the whole



- ▶ Need to exploit symmetry