

ENGR-E 511; ENGR-E 399

Machine Learning for Signal Processing

Module 08:

Nonlinear Methods

Minje Kim

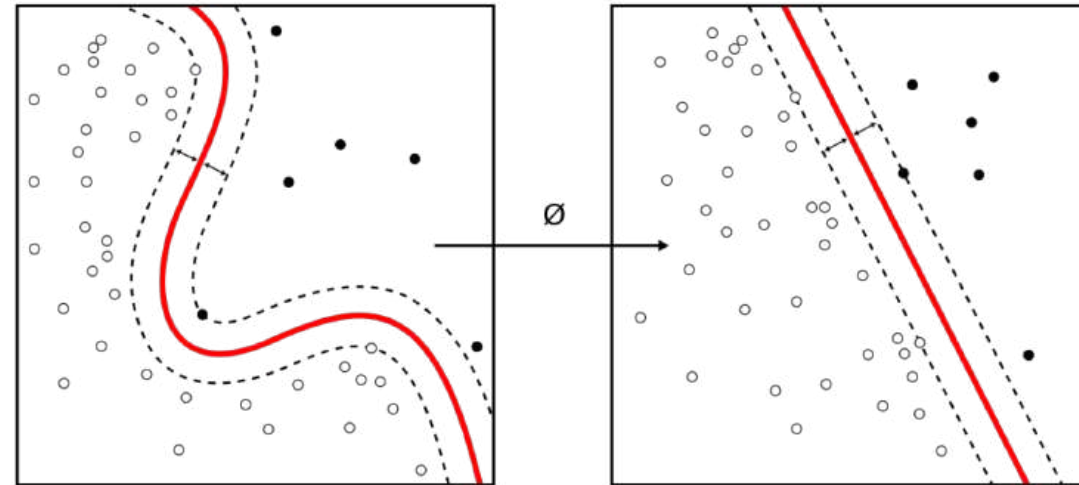
Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



https://en.wikipedia.org/wiki/Kernel_method

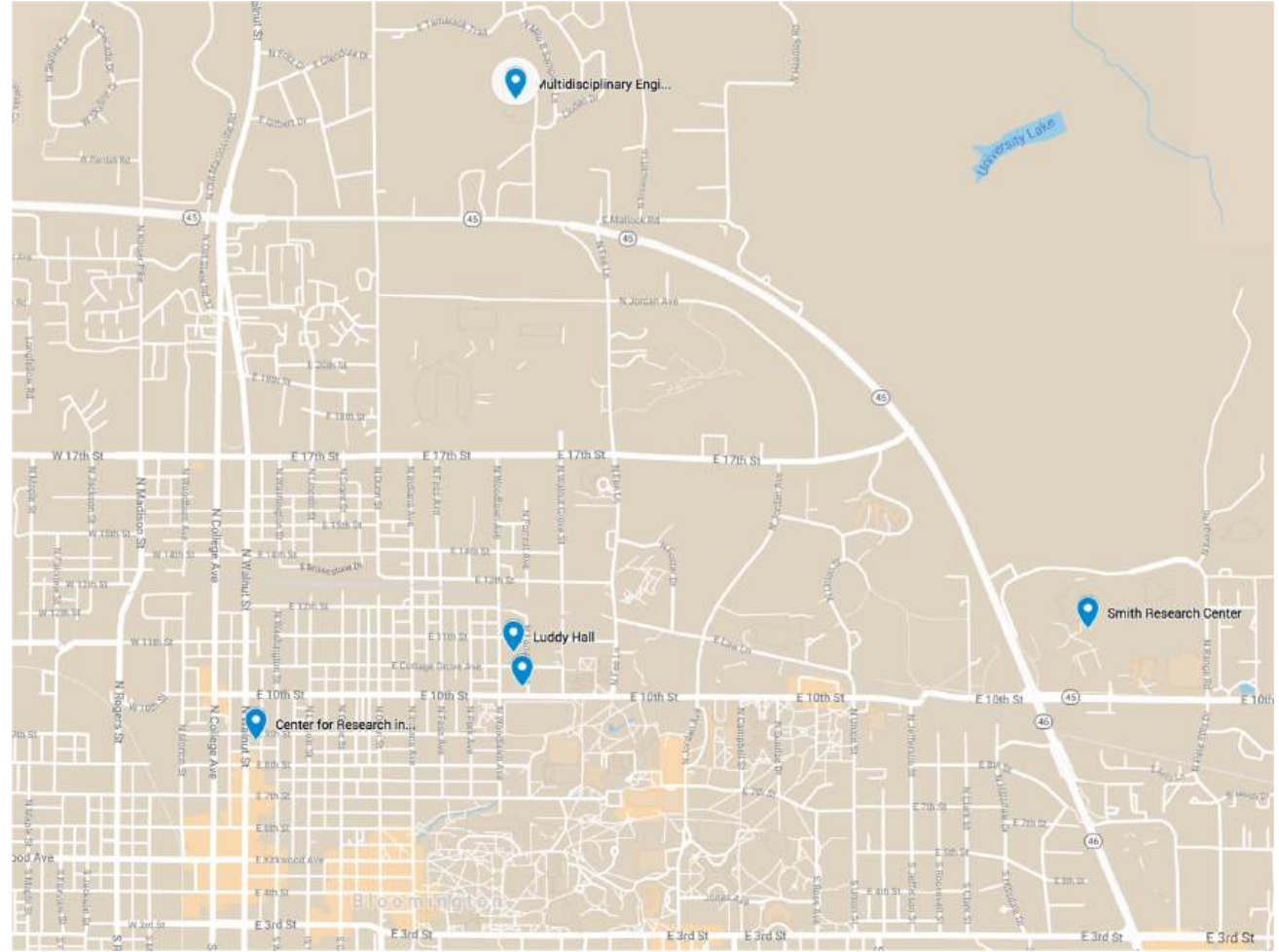
INDIANA UNIVERSITY
**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**



A Motivating (Linear) Example

- SICE is dispersed

- I'm finally moved to Luddy Hall
 - Yet, SICE is dispersed
- What if I don't know their locations but their pairwise distance?
 - Can I still recover the map?

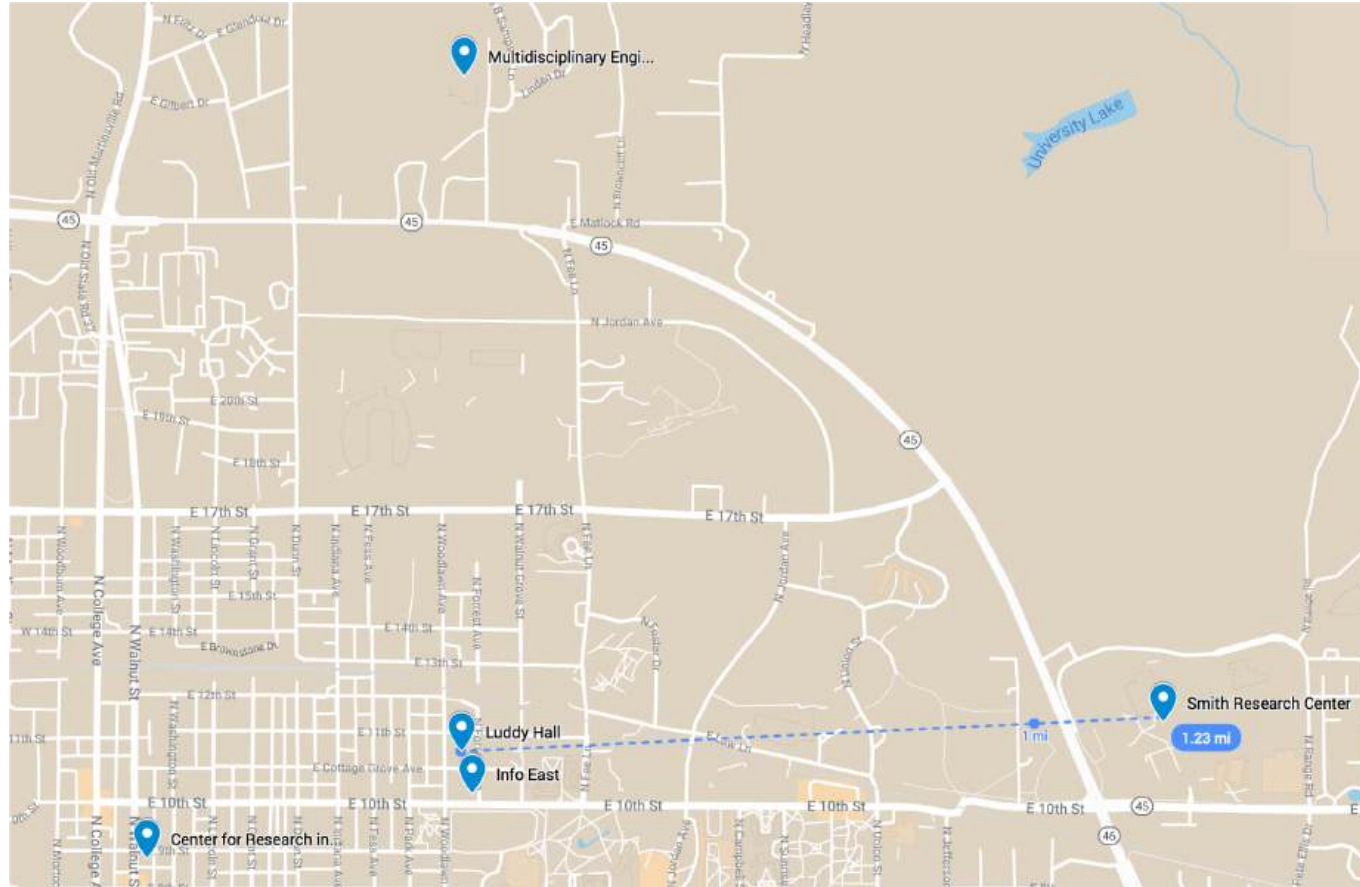


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

A Motivating (Linear) Example

- SICE is dispersed



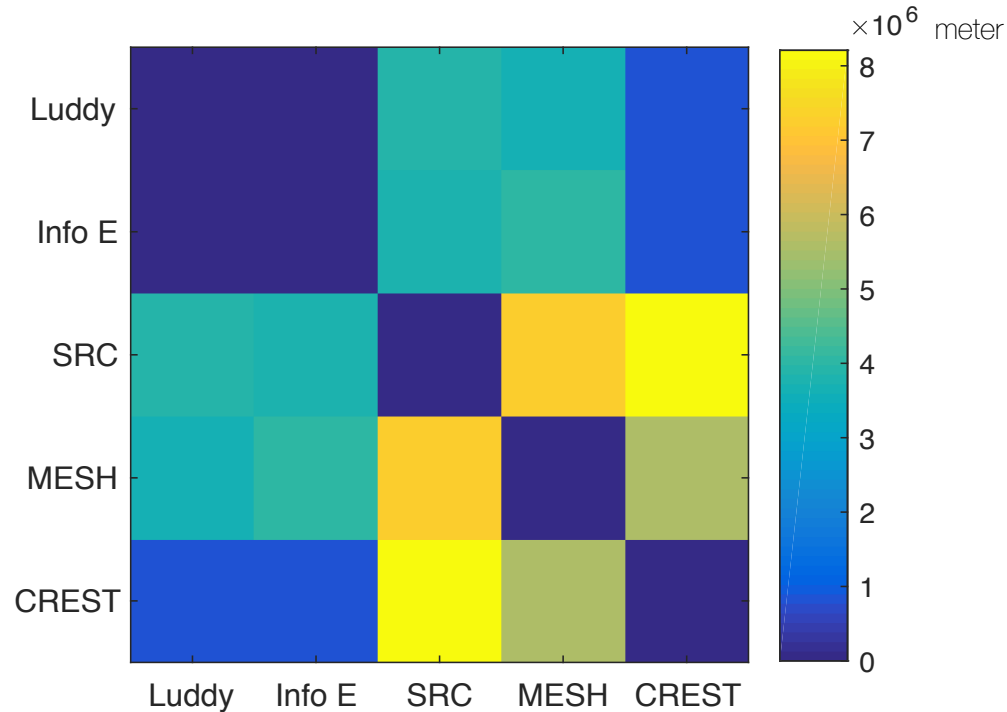
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

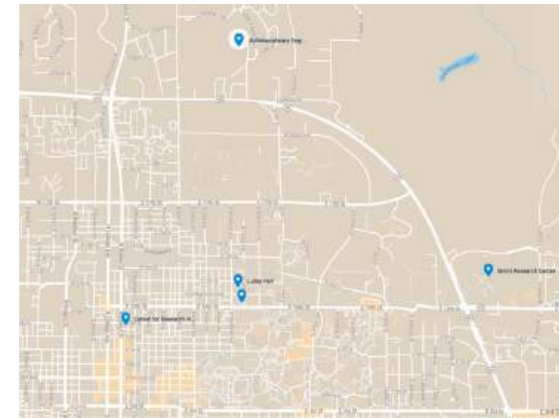
A Motivating (Linear) Example

- SICE is dispersed

- From now on I pretend like I don't know their locations
- Instead, I prepare a pairwise distance matrix



- From this, I want to recover their locations



MultiDimensional Scaling (MDS)

- A linear model

- The (squared) Euclidean pairwise distance matrix

$$M_{ij} = ||\mathbf{x}_i - \mathbf{x}_j||^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_i + \mathbf{x}_j^\top \mathbf{x}_j - 2\mathbf{x}_i^\top \mathbf{x}_j$$

- If it were not for the quadratic terms, we can say $\mathbf{W} \approx \mathbf{y}_i^\top \mathbf{y}_j$ ← Inner product between ESTIMATED locations (similarity)
- So, the goal of MDS is to transform M into W by eliminating $\mathbf{x}_i^\top \mathbf{x}_i + \mathbf{x}_j^\top \mathbf{x}_j$
- Suppose the average of M along the row

$$\bar{M} = \frac{1}{N} \sum_j M_{ij} = (N\mathbf{x}_i^\top \mathbf{x}_i + \sum_j \mathbf{x}_j^\top \mathbf{x}_j - 2\mathbf{x}_i^\top \sum_j \mathbf{x}_j) / N$$

- Then, let's subtract it from M

$$\tilde{M}_{ij} = M_{ij} - \bar{M}_i = \mathbf{x}_j^\top \mathbf{x}_j - \frac{1}{N} \sum_j \mathbf{x}_j^\top \mathbf{x}_j - 2\mathbf{x}_i^\top \mathbf{x}_j + 2\mathbf{x}_i^\top \bar{\mathbf{x}}$$

- Suppose the average of \tilde{M} along the column

$$|\tilde{M}_{ij} = \mathbf{x}_j^\top \mathbf{x}_j - \frac{1}{N} \sum_j \mathbf{x}_j^\top \mathbf{x}_j - 2\bar{\mathbf{x}}^\top \mathbf{x}_j + 2\bar{\mathbf{x}}^\top \bar{\mathbf{x}}$$

- Then, let's subtract it from \tilde{M}

$$\mathbf{W}_{ij} = \tilde{M}_{ij} - |\tilde{M}_j = -2\mathbf{x}_i^\top \mathbf{x}_j + 2\mathbf{x}_i^\top \bar{\mathbf{x}} + 2\bar{\mathbf{x}}^\top \mathbf{x}_j - 2\bar{\mathbf{x}}^\top \bar{\mathbf{x}} = -2(\mathbf{x}_i - \bar{\mathbf{x}})^\top (\mathbf{x}_j - \bar{\mathbf{x}})$$

MultiDimensional Scaling (MDS)

- A linear model

- We start from the distance matrix M
 - Subtract the average of rows and columns sequentially
- As a result we get a matrix W
 - Inner products of centered data samples (let's discard -2 from the previous slide)

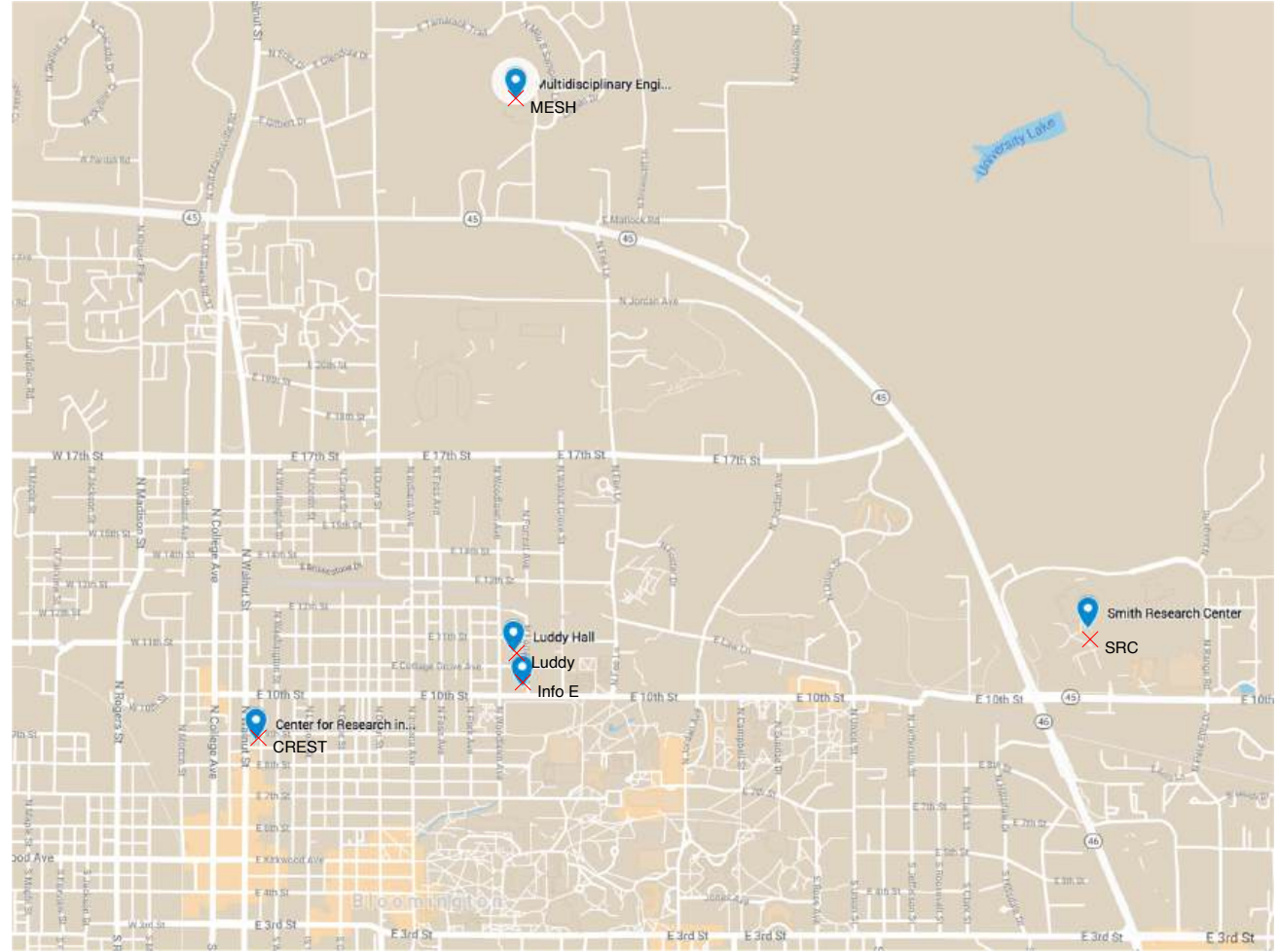
$$W_{ij} = (x_i - \bar{x})^\top (x_j - \bar{x}) \quad W = Y^\top Y \quad Y = X - \bar{x}$$

- So what?
 - W is a symmetric positive definite matrix
 - If we can decompose it into two factors, we're done. Any ideas?
 - Eigendecomposition! $W = UVU^\top = \left(UV^{\frac{1}{2}}\right) \left(V^{\frac{1}{2}}U^\top\right)$
- **The recipe for MDS**
 - Convert your pairwise distance matrix M into a positive definite matrix W
 - Do eigendecomposition on the matrix W
 - First few vectors of $UV^{\frac{1}{2}}$ holds recovers the locations (coordinates)

MDS on SICE Buildings

- Using two eigenvectors

- Note that there are rotation and shift ambiguity in MDS solutions



INDIANA UNIVERSITY

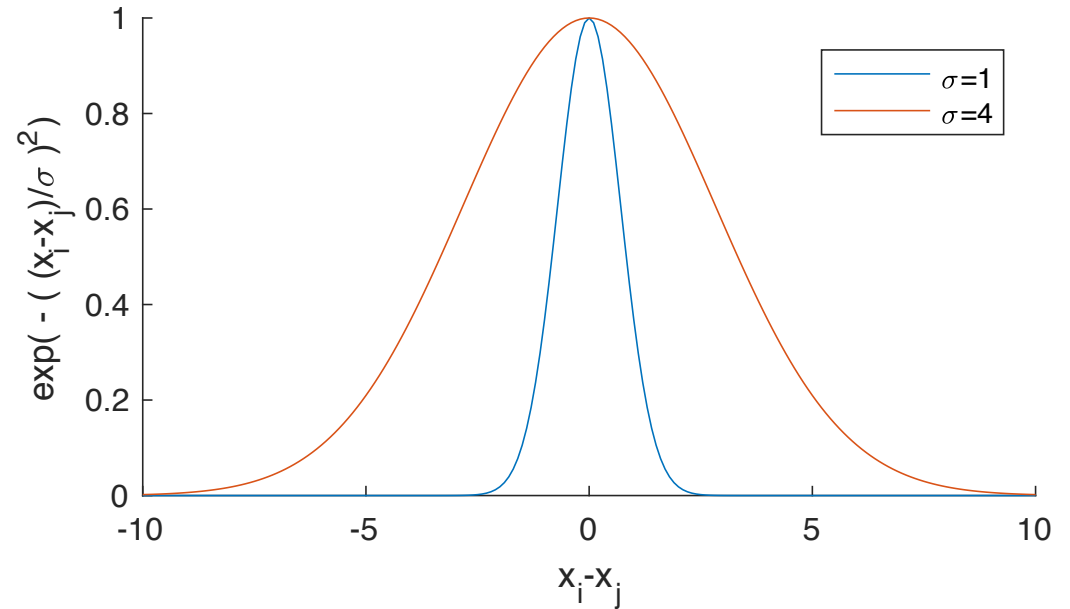
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Euclidean Distance?

- MDS and beyond

- Inner products can define some kind of affinity
- How about the other types of affinity?
- A popular choice:

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^2}\right)$$

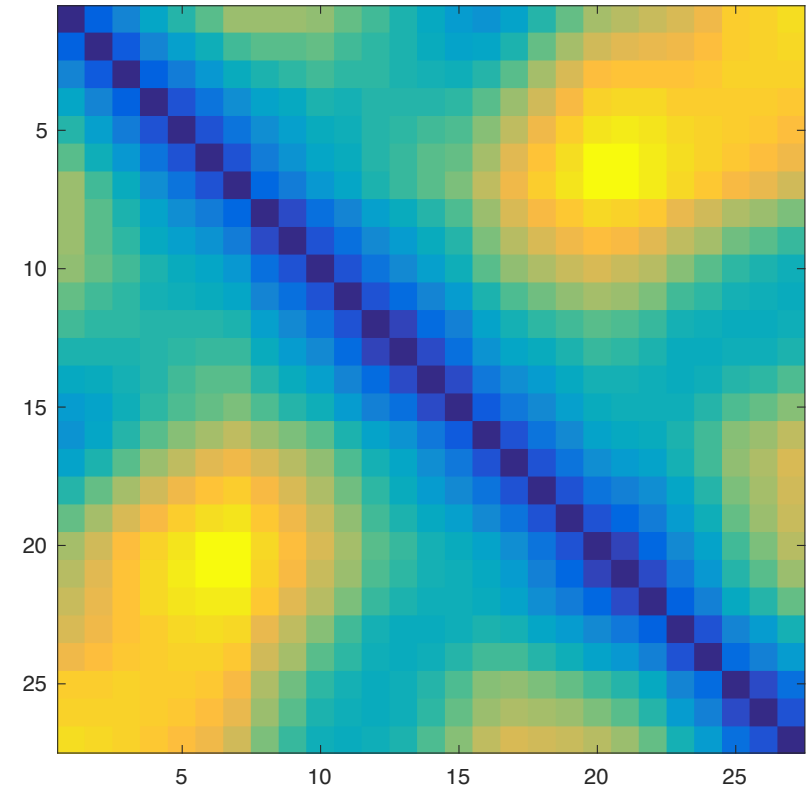
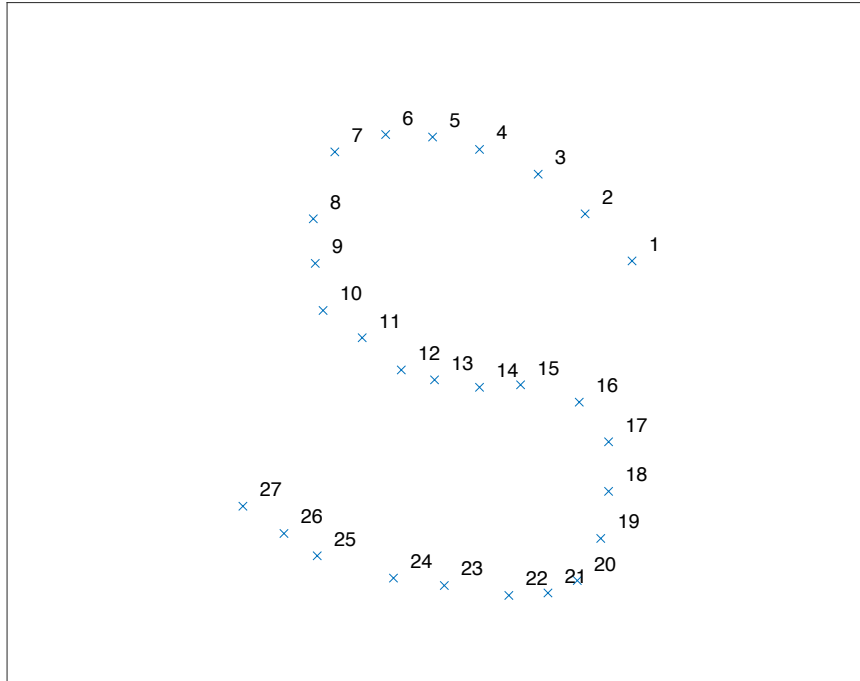


- If σ is small, closer data points become more similar
 - It weighs more on the closer relationships while suppressing farther ones
 - We call it **Radial Basis Function (RBF)**

Euclidean Distance?

- MDS and beyond

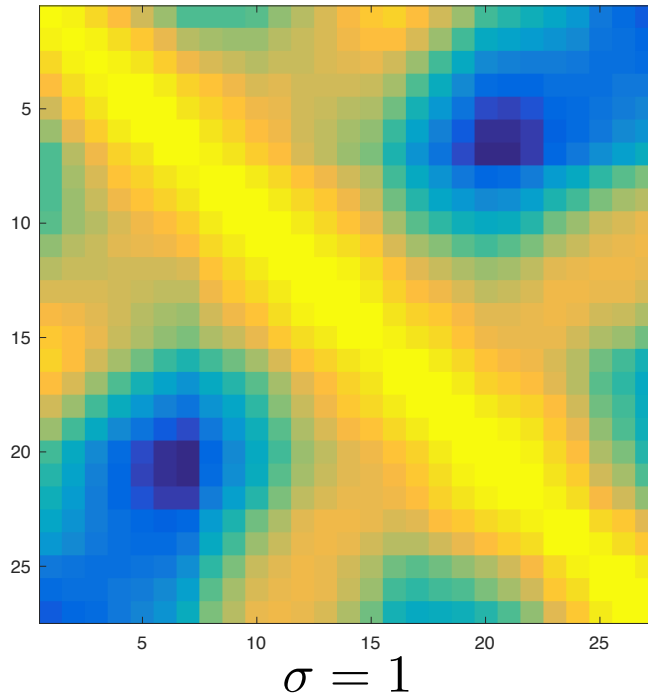
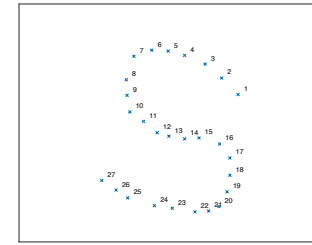
- An “S” shape and the pairwise Euclidean distance matrix



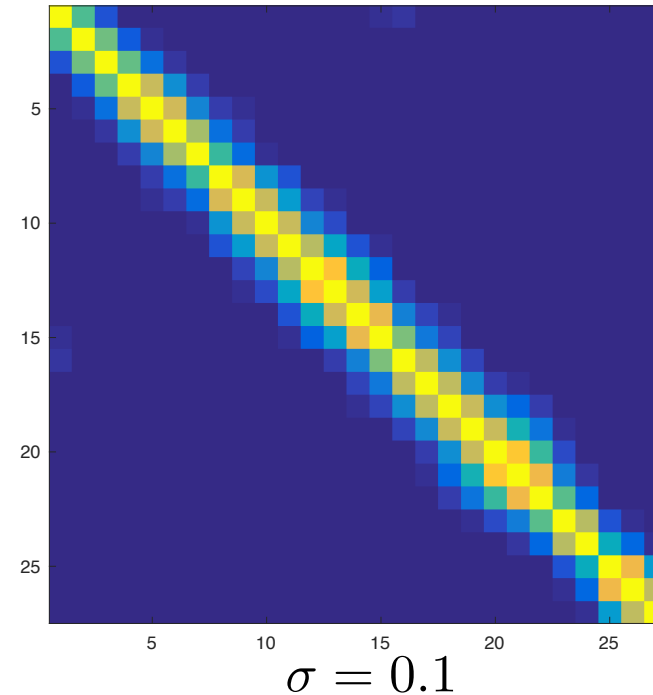
Euclidean Distance?

- MDS and beyond

- The affinity matrix with RBF kernels



$\sigma = 1$



$\sigma = 0.1$

- The narrower kernel preserves the local affinity better
 - While largely ignoring the global (linear) structure

Nonlinear Similarity

- When do we care about this?

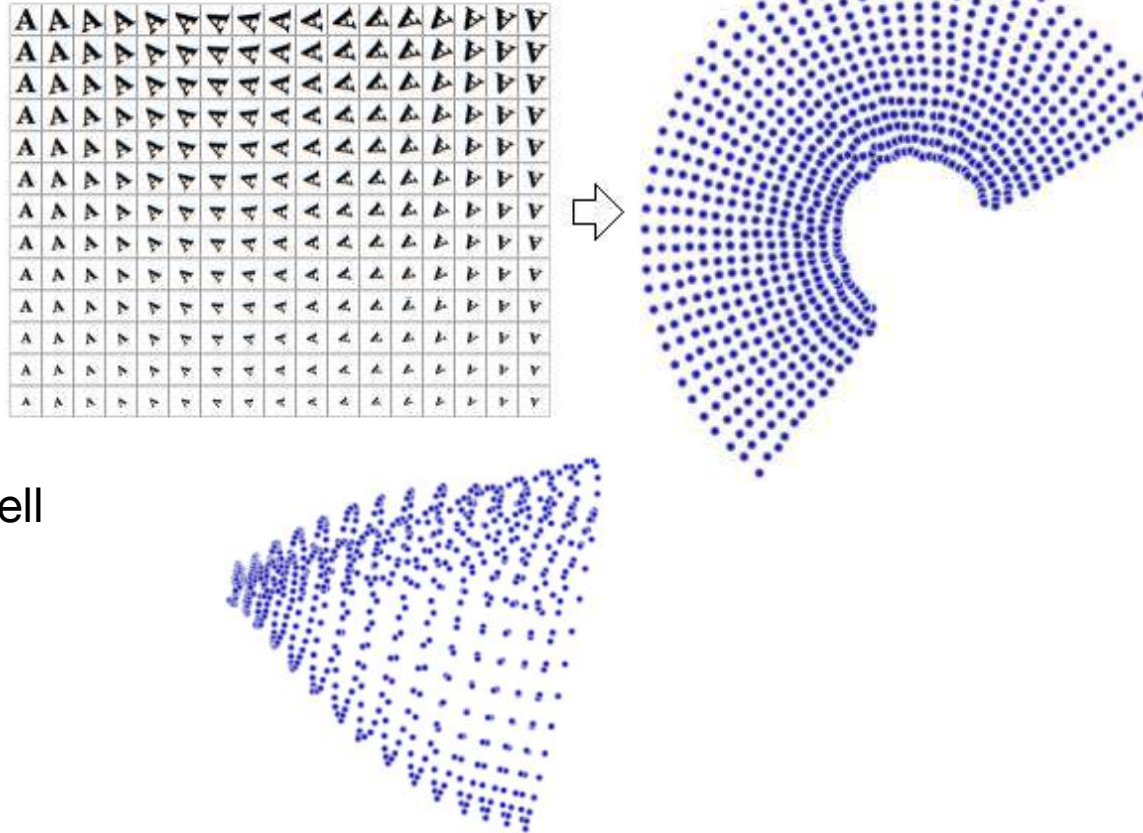
- MDS can recover the map from a kind of affinity matrix
 - We can control the dimensionality of the map using eigenvalues
 - e.g. The buildings are actually lying in a 3D space, not 2D
 - MDS is a linear technique and its affinity matrix is inner product
 - Perfect reconstruction is possible if we know the rank
- From this, what we can say is
 - We can recover the lower dim representation from the affinity matrix
 - Not from the data points
 - We can use whatever funky affinity matrix
 - If we don't care about the perfect reconstruction
- So, when do we not care about the perfect reconstruction?
 - If the data points are too noisy and high dimensional
 - Noisy AND high dimensional?



Intrinsic Dimensionality

- Moving A's

- How many dimensions do we need?

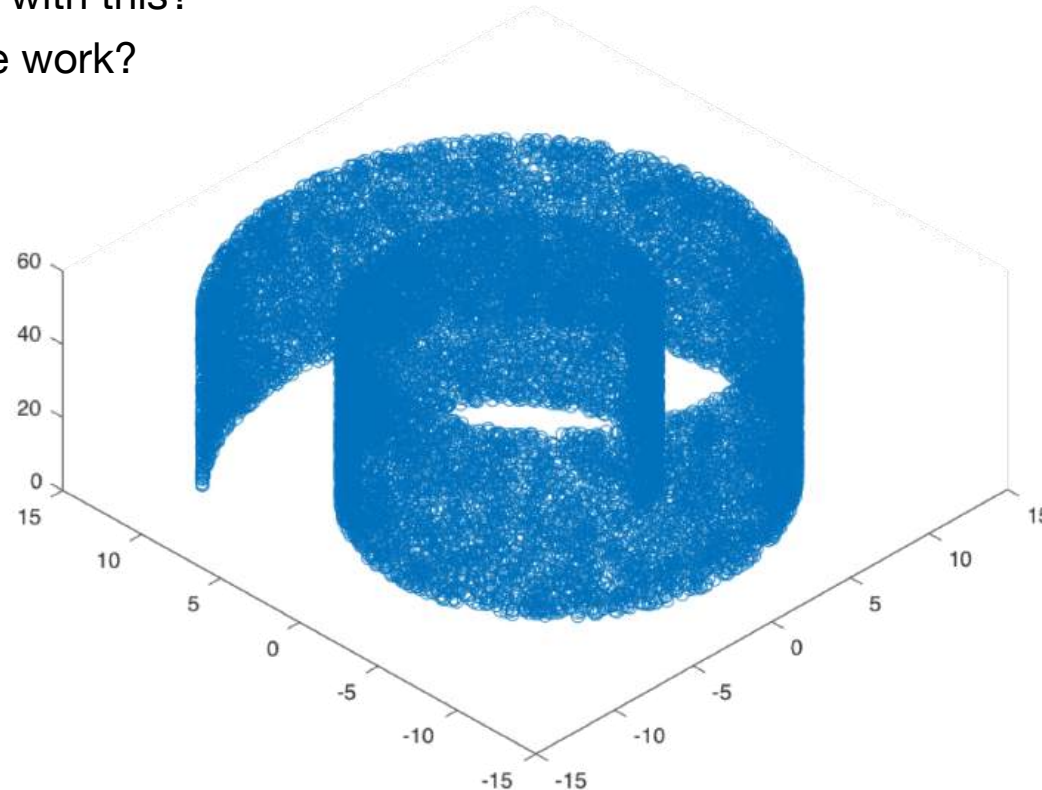


- PCA doesn't do it well

Manifold Learning

- Nonlinear dimension reduction: the Swiss roll example

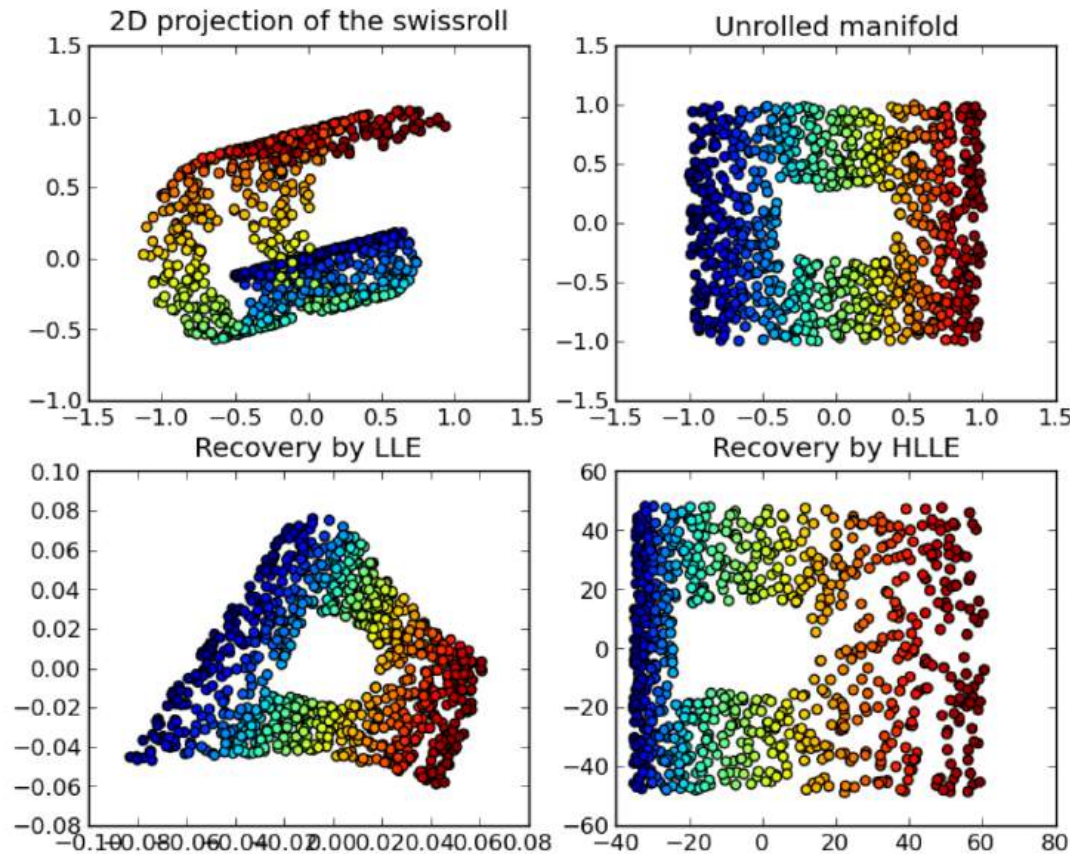
- The Swiss roll example
 - What do you want to do with this?
 - Would a linear technique work?



Manifold Learning

- Nonlinear dimension reduction: the Swiss roll example

○ Unrolling needs a nonlinear projection



Laplacian Eigenmap

- Eigendecomposition on the RBF kernel

- First, we build an affinity matrix using the RBF kernel

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^2} \right)$$



- With a small σ for the locality
- In practice, we skip the non-neighboring pairs $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = 0$
 - If $\|\mathbf{x}_i - \mathbf{x}_j\|^2 > \tau$
 - Or, if $\mathbf{x}_j \notin \mathcal{N}_i$
- The objective function $\arg \min_{\mathbf{Y}} \sum_{i,j} \frac{(\mathbf{y}_i - \mathbf{y}_j)^\top (\mathbf{y}_i - \mathbf{y}_j)}{\text{Distance after projection}} \frac{W_{ij}}{\text{Original affinity } \mathcal{K}}$
 - It will be minimal if originally similar pairs are also similar after the projection
- Why not the direct decomposition on \mathbf{W} (as in MDS)?
 - Because perfect reconstruction might not be possible

Laplacian Eigenmap

- Eigendecomposition on the RBF kernel

○ The objective function $\arg \min_{\mathbf{Y}} \sum_{i,j} (\mathbf{y}_i - \mathbf{y}_j)^\top (\mathbf{y}_i - \mathbf{y}_j) \mathbf{W}_{ij}$

○ Optimization $\mathcal{E} = \sum_{i,j} (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_j^\top \mathbf{y}_i + \mathbf{y}_j^\top \mathbf{y}_j) \mathbf{W}_{ij}$
 $= \sum_i \mathbf{y}_i^\top \mathbf{y}_i \sum_j \mathbf{W}_{ij} + \sum_j \mathbf{y}_j^\top \mathbf{y}_j \sum_i \mathbf{W}_{ij} - 2 \sum_{i,j} \mathbf{y}_j^\top \mathbf{y}_i \mathbf{W}_{ij}$
 $= 2(\mathbf{Y} \mathbf{D} \mathbf{Y}^\top - \mathbf{Y} \mathbf{W} \mathbf{Y}^\top) = 2\mathbf{Y} \mathbf{L} \mathbf{Y}^\top$

$\mathbf{D}_{ii} = \sum_j \mathbf{W}_{ij}$
A diagonal matrix  
A vector

$$\mathbf{L} = \mathbf{D} - \mathbf{W}$$

○ To avoid trivial solution $\mathbf{Y} = \mathbf{0}$ we specify the norm of the solution

$$\arg \min_{\mathbf{Y}} \mathbf{Y} \mathbf{L} \mathbf{Y}^\top \quad \text{s.t.} \quad \mathbf{Y} \mathbf{D} \mathbf{Y}^\top = \mathbf{I}$$

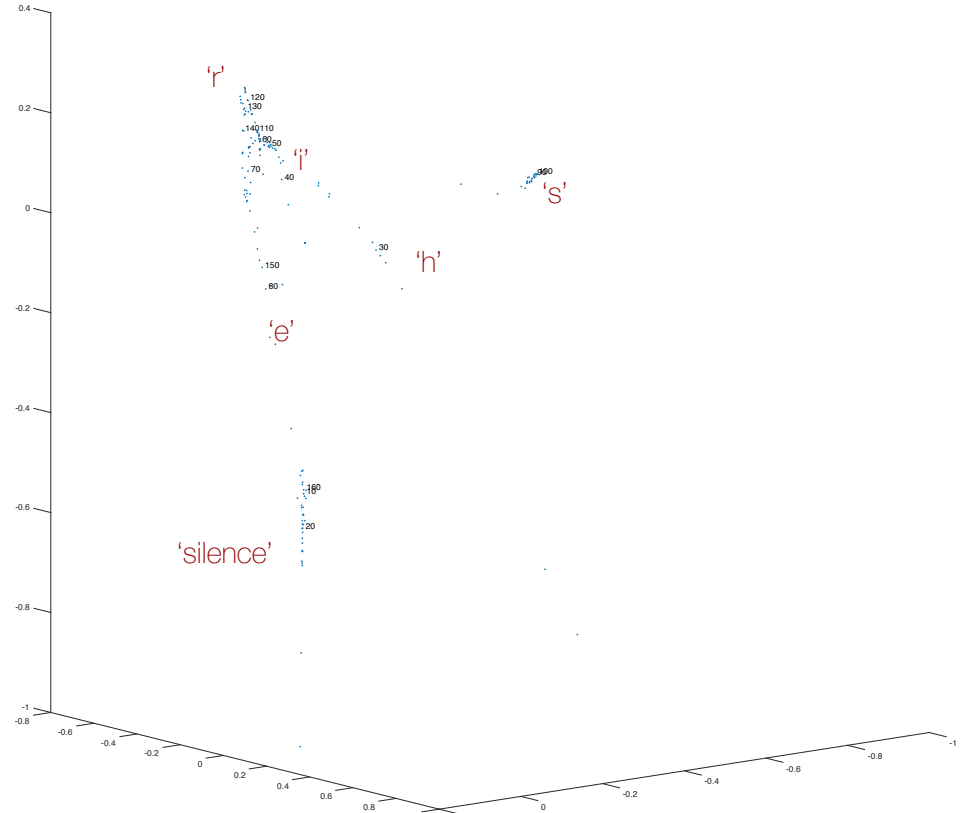
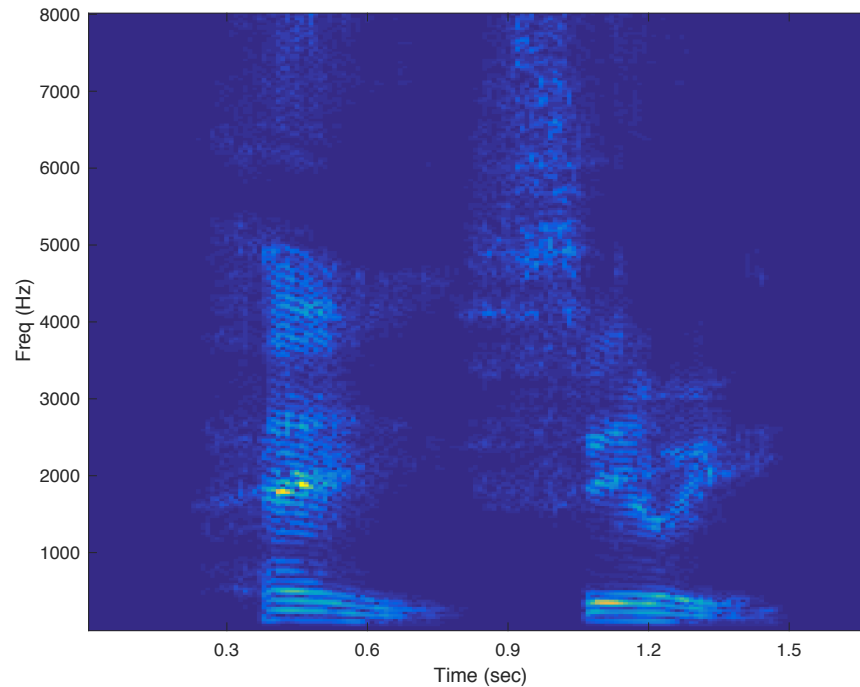
○ Or, $\arg \min_{\tilde{\mathbf{Y}}} \tilde{\mathbf{Y}} \tilde{\mathbf{L}} \tilde{\mathbf{Y}}^\top \quad \text{s.t.} \quad \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top = \mathbf{I} \quad \tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}} \quad \tilde{\mathbf{Y}} = \mathbf{Y} \mathbf{D}^{\frac{1}{2}}$

□ Ring a bell?

- Eigendecomposition on the Laplacian matrix $\tilde{\mathbf{L}}$

Laplacian Eigenmaps

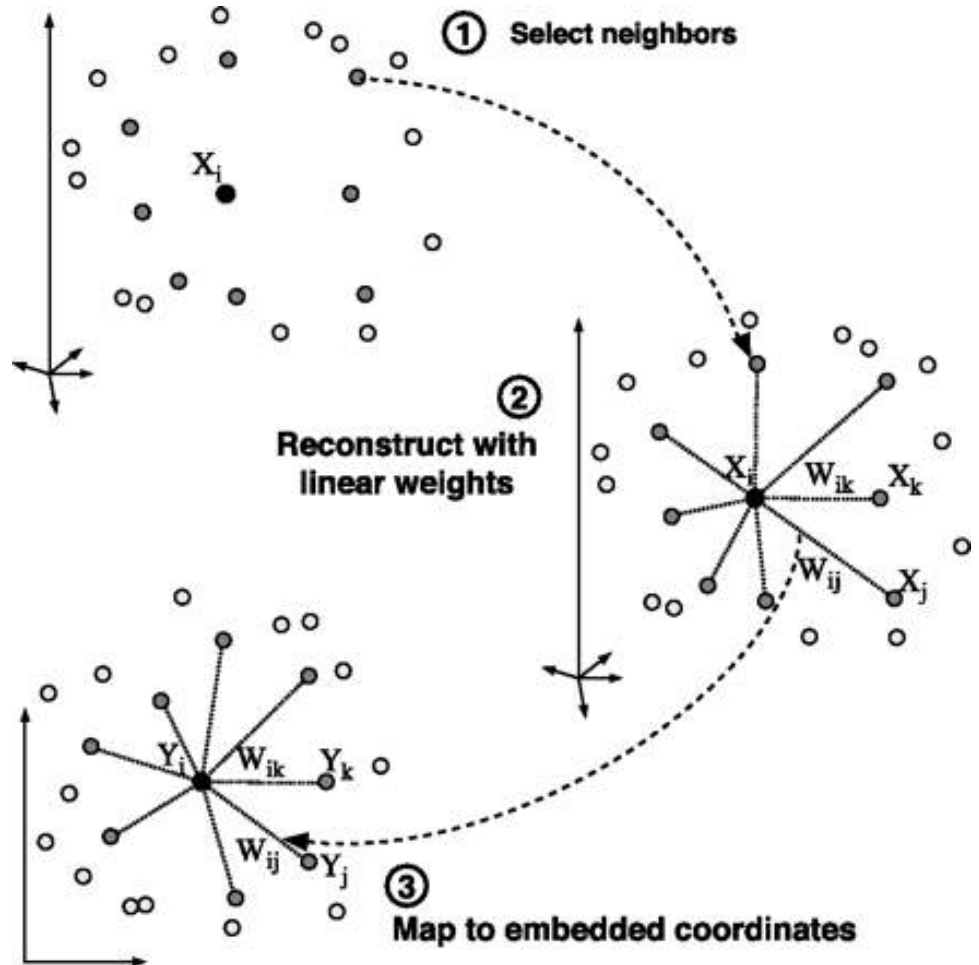
- On “Hey, Siri”



Locally Linear Embedding (LLE)

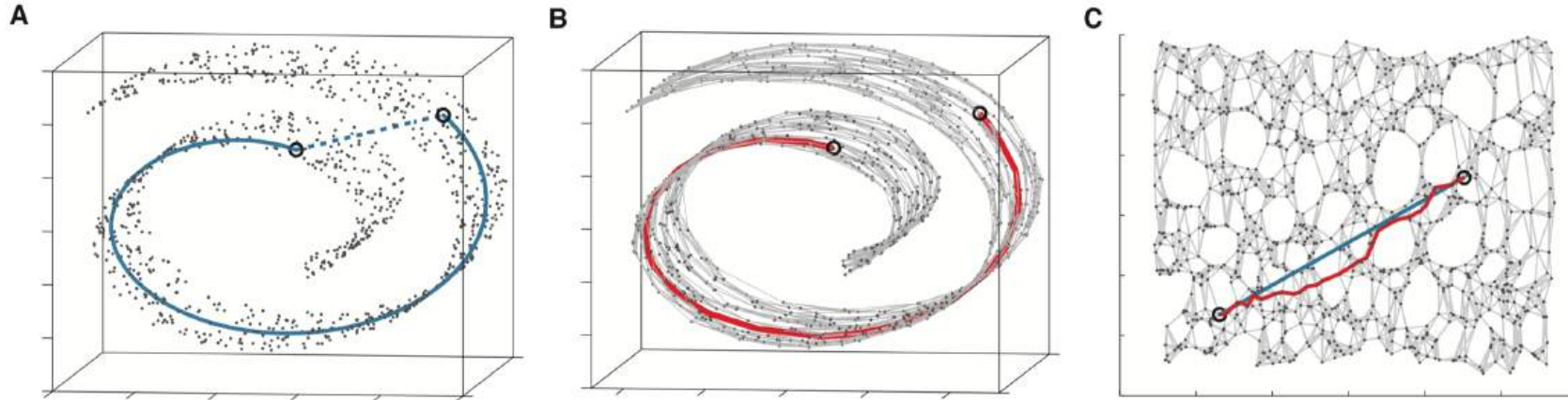
- The other manifold learning techniques

- For each data points, find a set of nearest neighbors
- Construct a local convex combination
- Find the lower dimensional version that preserves these local structures
- Yet another eigendecomposition problem



ISOMAP

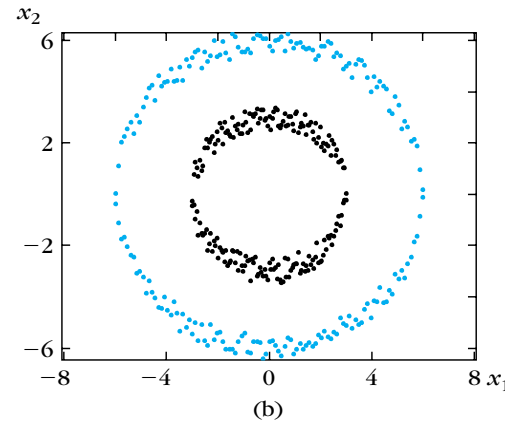
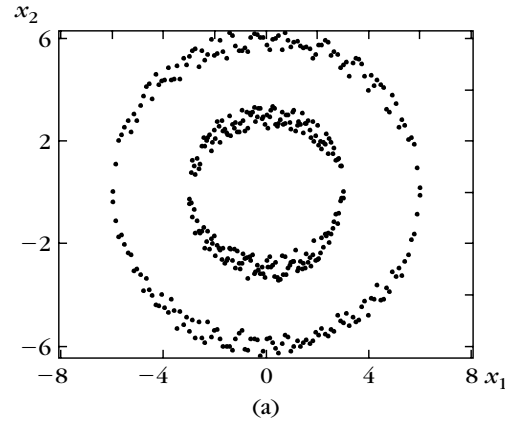
- The other manifold learning techniques



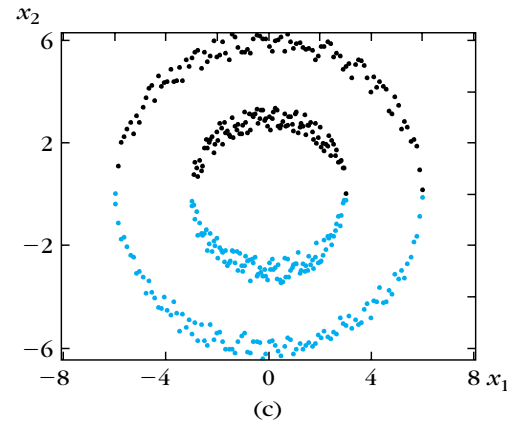
- For each data points, find a set of nearest neighbors
- Compute pairwise shortest paths (figure B)
 - Using Dijkstra's algorithm
 - Which approximates the geodesic distance (figure A)
- You built a nonlinear distance matrix. Do MDS on this.

Spectral Clustering

- Nonlinear clustering



Spectral Clustering



kMeans

Eigendecomposition on the Kernel Matrix

- What is it good for?

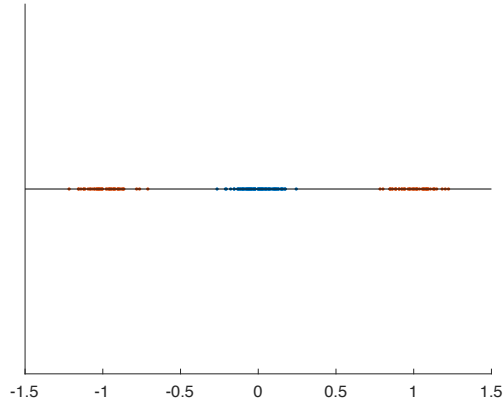
- The techniques we learned:
 - We can find a lower dimensional feature space from the affinity matrix, not from the data
 - Eigenvectors correspond to the coefficients after the projection, not the projection themselves
 - The projection can be nonlinear, too
 - And we don't have an access to the projection
- What's so special about those kernel matrices?
 - We only saw RBF kernel, but there are other options
- In general...
 - If we use the kernel tricks we can avoid learning direct nonlinear transforms
 - Kernel tricks still let us discover nonlinear structure from the data
 - We call this kind of algorithms **kernel methods**



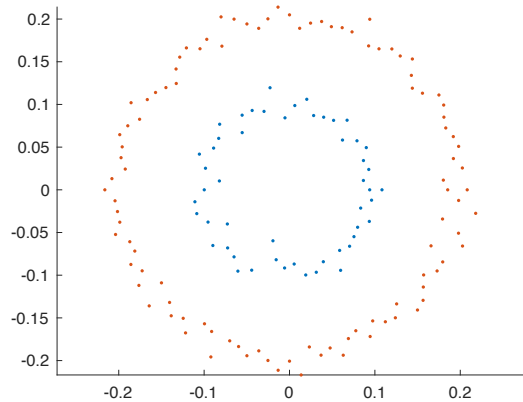
Nonlinear Transformation

- For linear separability

○ A problem not linearly separable

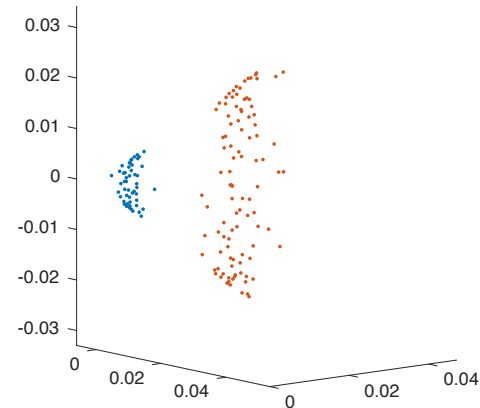
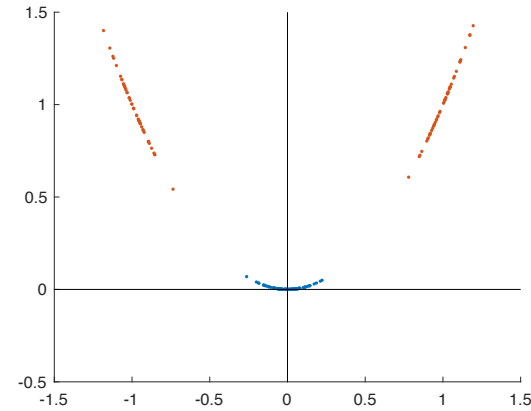


$$\phi(x_1) = [x_1, x_1^2]^\top$$



$$\phi(x_1, x_2) = [x_1^2, x_2^2, x_1x_2]^\top$$

○ After the transform



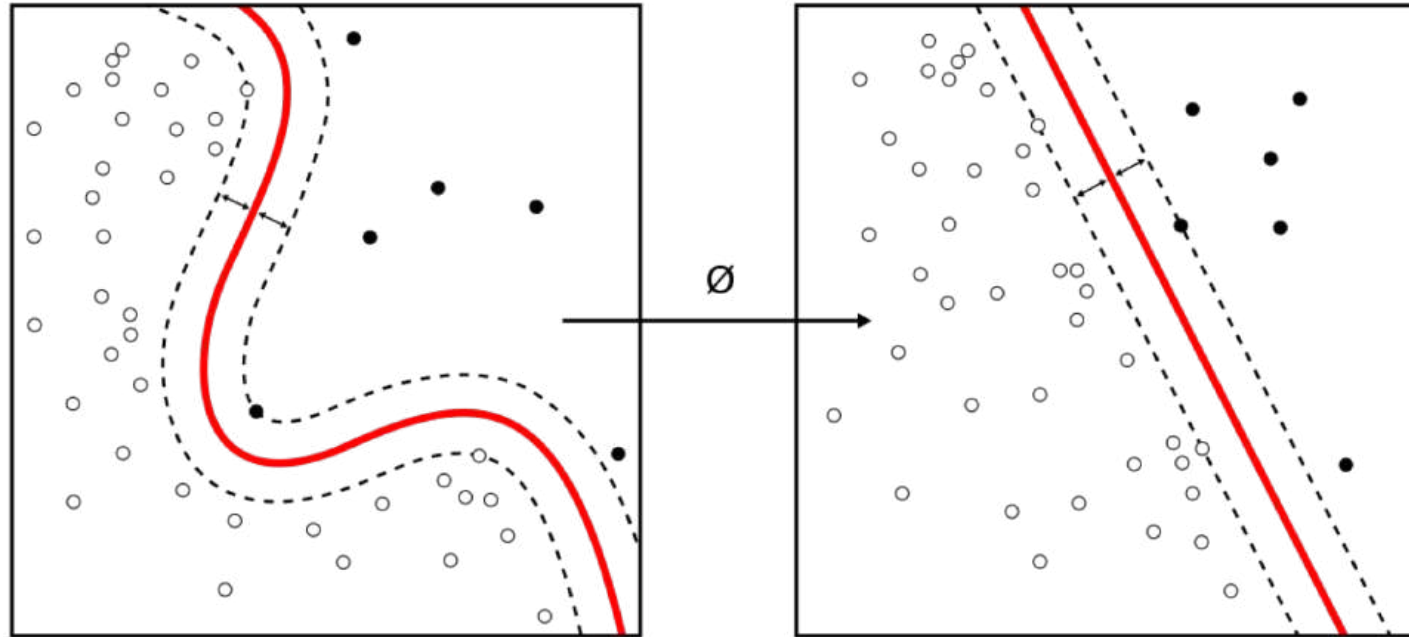
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Nonlinear Transformation

- For linear separability

- Nonlinearly separable functions might be able to be solved linearly
 - In the nonlinearly found feature space



- If you increase the dimensionality during the nonlinear transform, you hope that you can solve it linearly

Nonlinear Transformation

- For linear separability

- It's difficult to find out this magical nonlinear transform function
 - Because it's problem-specific
- What we know is that the problem becomes easier to solve if we blow up the dimensionality
- Then, apply an arbitrary transform function to increase the dimensionality?
 - I don't want to do that
 - It may make the problem complicated with the high dimensionality
- Is there any way to work with the high dimensional data without explicitly knowing their values?
 - Kernels!
 - The kernel matrix can ***imply*** the relationship in the ***transformed feature space***




Radian Basis Function (RBF)

- An infinite dimensionality

- We've seen the use of RBF kernels in the previous examples
 - What's so special about it? (you know, we still don't know the mapping function)

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^2} \right)$$

- The dimensionality of the feature space of the RBF kernel is **infinite**
 - So, using the RBF kernel implies that the transform increased the dimensionality a lot

- In other words $\mathbf{z}_i = \phi(\mathbf{x}_i)$ $\mathbf{z}_j = \phi(\mathbf{x}_j)$  We don't explicitly know the mapping

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{z}_i^\top \mathbf{z}_j$$
  Kernel is the inner products of the features

$$\mathbf{z}_i, \mathbf{z}_j \in \mathbb{R}^\infty$$
  For RBF the feature space is with an infinite dimensionality

Radian Basis Function (RBF)

- An infinite dimensionality

○ Let's prove it for a simple case

□ x_i is scalar

□ $\sigma = 1$

○ Then the kernel matrix can be defined as follows

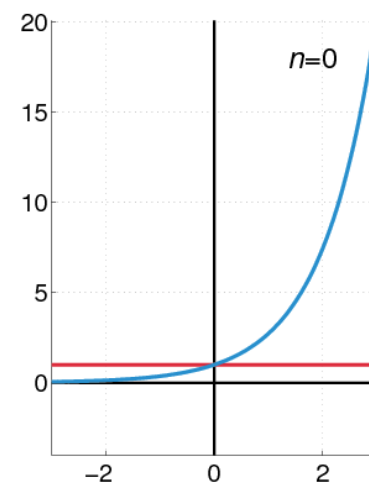
$$\begin{aligned}\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) &= \exp(-(x_i - x_j)^2) = \exp(-x_i^2 - x_j^2 + 2x_i x_j) \\ &= \exp(-x_i^2) \exp(-x_j^2) \exp(2x_i x_j)\end{aligned}$$

Taylor series $f(x) \approx f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \frac{f'''(a)}{3!}(x - a)^3 + \dots$

$$= \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n$$

Therefore $\exp(x) \approx \sum_{n=0}^{\infty} \frac{\exp(0)}{n!} (x - 0)^n = \sum_{n=0}^{\infty} \frac{x^n}{n!}$

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-x_i^2) \exp(-x_j^2) \sum_{n=0}^{\infty} \frac{2^n x_i^n x_j^n}{n!}$$



Radian Basis Function (RBF)

- An infinite dimensionality

- An inner product between infinite-dimensional vectors

$$\sum_{n=0}^{\infty} \frac{2^n x_i^n x_j^n}{n!} = \left[\frac{\sqrt{(2^0)}x_i^0}{\sqrt{0!}}, \frac{\sqrt{(2^1)}x_i^1}{\sqrt{1!}}, \frac{\sqrt{(2^2)}x_i^2}{\sqrt{2!}}, \dots, \frac{\sqrt{(2^\infty)}x_i^\infty}{\sqrt{\infty!}} \right] \begin{bmatrix} \frac{\sqrt{(2^0)}x_j^0}{\sqrt{0!}} \\ \frac{\sqrt{(2^1)}x_j^1}{\sqrt{1!}} \\ \frac{\sqrt{(2^2)}x_j^2}{\sqrt{2!}} \\ \dots \\ \frac{\sqrt{(2^\infty)}x_j^\infty}{\sqrt{\infty!}} \end{bmatrix}$$

- Eventually, the RBF kernel is an inner product between infinite dimensional features

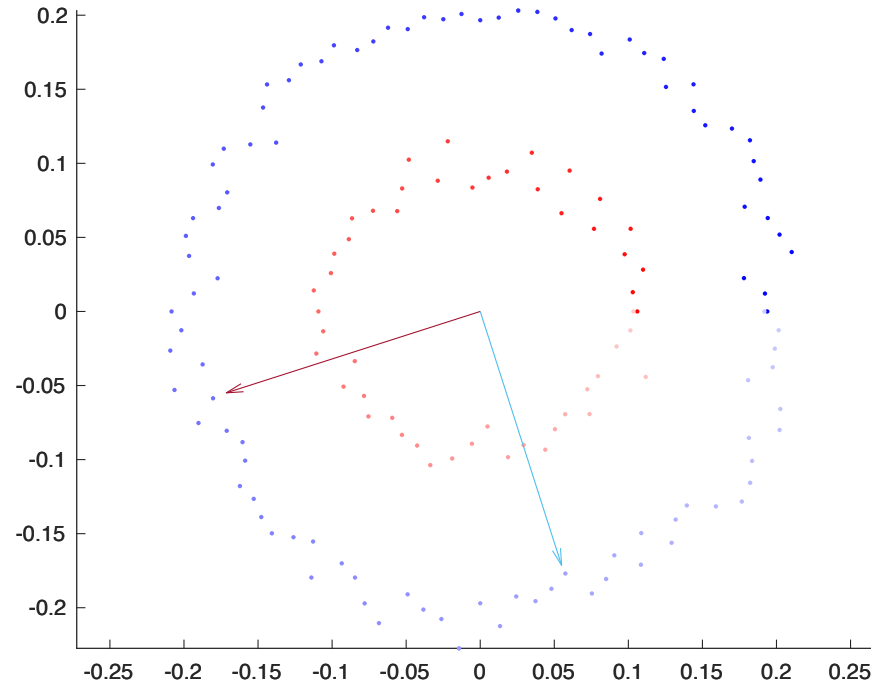
$$\mathcal{K}(x_i, x_j) = \exp(-x_i^2) \exp(-x_j^2) \sum_{n=0}^{\infty} \frac{2^n x_i^n x_j^n}{n!}$$

$$\mathcal{K}(x_i, x_j) = \left[\exp(-x_i^2) \frac{\sqrt{(2^0)}x_i^0}{\sqrt{0!}}, \exp(-x_i^2) \frac{\sqrt{(2^1)}x_i^1}{\sqrt{1!}}, \exp(-x_i^2) \frac{\sqrt{(2^2)}x_i^2}{\sqrt{2!}}, \dots, \exp(-x_i^2) \frac{\sqrt{(2^\infty)}x_i^\infty}{\sqrt{\infty!}} \right] \begin{bmatrix} \exp(-x_j^2) \frac{\sqrt{(2^0)}x_j^0}{\sqrt{0!}} \\ \exp(-x_j^2) \frac{\sqrt{(2^1)}x_j^1}{\sqrt{1!}} \\ \exp(-x_j^2) \frac{\sqrt{(2^2)}x_j^2}{\sqrt{2!}} \\ \dots \\ \exp(-x_j^2) \frac{\sqrt{(2^\infty)}x_j^\infty}{\sqrt{\infty!}} \end{bmatrix}$$

Kernel PCA

- Concentric circles

- What would be the principal components of these data samples?

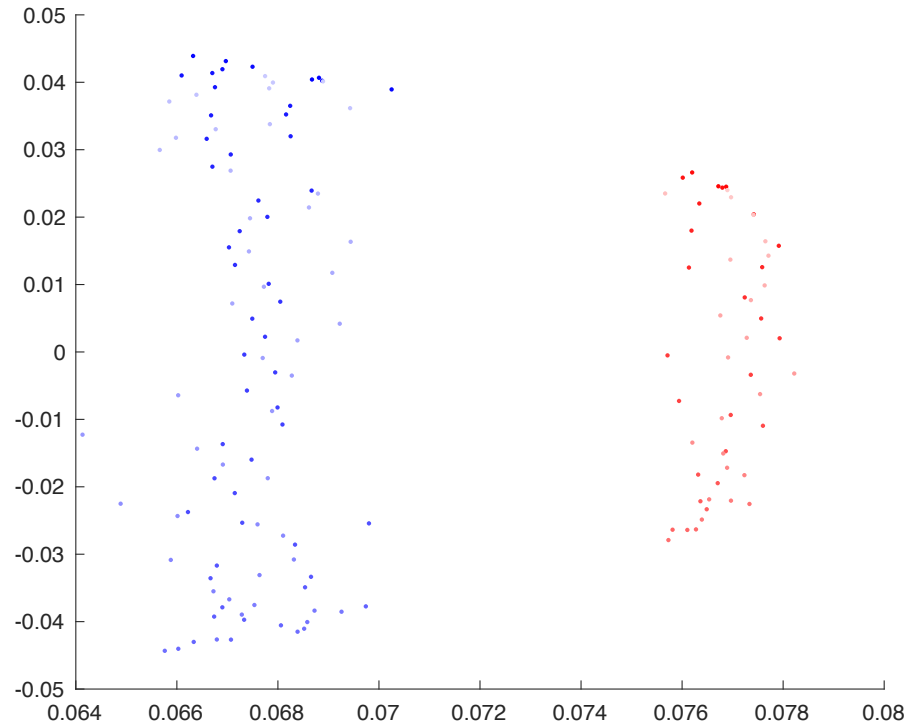


- An ordinary PCA doesn't give a satisfactory result

Kernel PCA

- Concentric circles

- On the other hand, kernel PCA can find out the (nonlinear) principal components



Kernel PCA

- Eigendecomposition with the kernel trick

- Suppose that we observe a dataset with a nonlinear structure
- Suppose that there is a nonlinear transform that corresponds to it

$$\mathbf{z}_i = \phi(\mathbf{x}_i)$$

- Although we don't know what it exactly would be
- We assume that \mathbf{z}_i is higher dimensional
 - And we assume the nonlinear structure turned into a linear structure
- Now we want to do PCA on \mathbf{z}_i , not on \mathbf{x}_i
 - Because \mathbf{z}_i is with a linear structure now, although high dimensional
- PCA on the nonlinear transformed data starts from a cov matrix

$$\mathbf{C} = \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^\top = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top$$

- But we don't exactly know the mapping function $\phi(\mathbf{x}_i)$
 - What can we do?

Kernel PCA

- Eigendecomposition with the kernel trick

- The definition of eigendecomposition gives us a hint

$$\begin{aligned} C\mathbf{v} &= \lambda\mathbf{v} && \text{What's their dimensionality?} \\ \left\{ \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i)\phi(\mathbf{x}_i)^\top \right\} \mathbf{v} &= \lambda\mathbf{v} \\ \frac{1}{N} \mathbf{Z}\mathbf{Z}^\top \mathbf{v} &= \lambda\mathbf{v} && \leftarrow \text{This is high dimensional, but a vector} \\ \mathbf{Z}\boldsymbol{\alpha} &= \mathbf{v} && \leftarrow \text{Plug in} \end{aligned}$$

- Eventually,

$$\begin{aligned} \frac{1}{N} \mathbf{Z}\mathbf{Z}^\top \mathbf{Z}\boldsymbol{\alpha} &= \lambda\mathbf{Z}\boldsymbol{\alpha} \\ \mathbf{Z}^\top \mathbf{Z}\boldsymbol{\alpha} &= N\lambda\boldsymbol{\alpha} \end{aligned}$$

The kernel matrix, not a covariance matrix

- We have a new eigendecomposition problem, on the Kernel matrix of the nonlinear transformed data

Kernel PCA

- Eigendecomposition with the kernel trick

- How do we recover the coefficients after the projection?
 - We want to do $y = v^\top Z$, but we know neither Z nor v
- With the kernel trick, we can do this, because $y = v^\top Z = \alpha^\top Z^\top Z = \alpha^\top \mathcal{K}$
- But we never know the actual projection vector $Z\alpha = v$
- Another thing to note
 - The eigendecomposition on the kernel matrix $\mathcal{K} = Z^\top Z$ will yield orthonormal eigenvectors α
 - But it doesn't guarantee the normalization of v
 - So, we need to normalize α further to make sure this: $v^\top v = \alpha^\top Z^\top Z \alpha = 1$
- One more thing to note: since the transform is nonlinear, the high dimensional features are not guaranteed to be centered

Kernel PCA

- The recipe

- Choose your kernel function
 - e.g. RBF
- Calculate the kernel matrix \mathcal{K}
 - Using the kernel function, not from the nonlinear transform
 - Center the kernel matrix
- Do eigendecomposition on \mathcal{K}
- Choose a few eigenvectors α
 - This will determine the dimensionality of this nonlinear dimensionality reduction algorithm
 - Normalize them further
- Compute the coefficients after the nonlinear dim reduction $y = \alpha^\top \mathcal{K}$
- **Note that this is exactly like MDS except we construct the kernel from the RBF function, not from an inner product**

Summary of the Kernel Methods

- Pros and cons

○ Pros

- We can do a nonlinear analysis of a data set
- We don't have to know the nonlinear feature transform

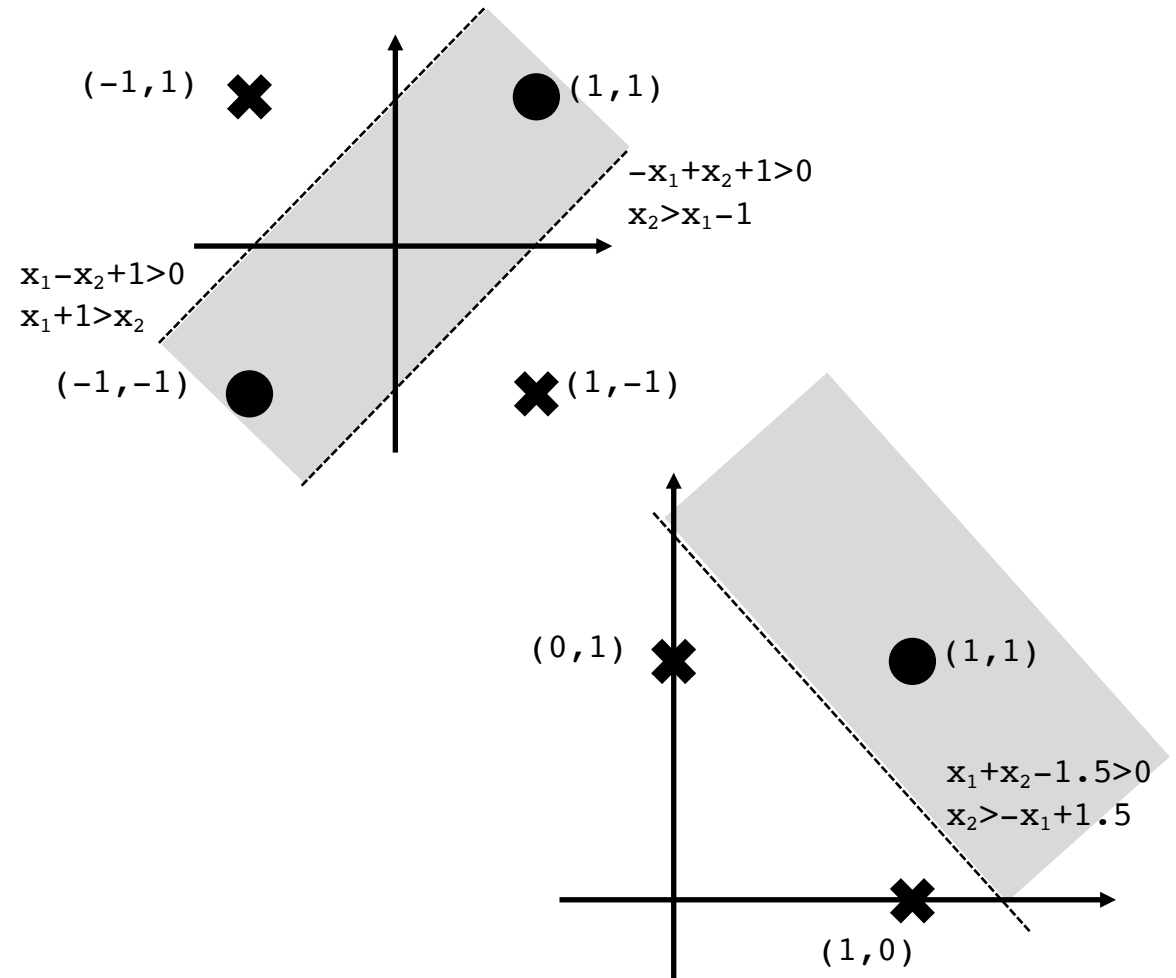
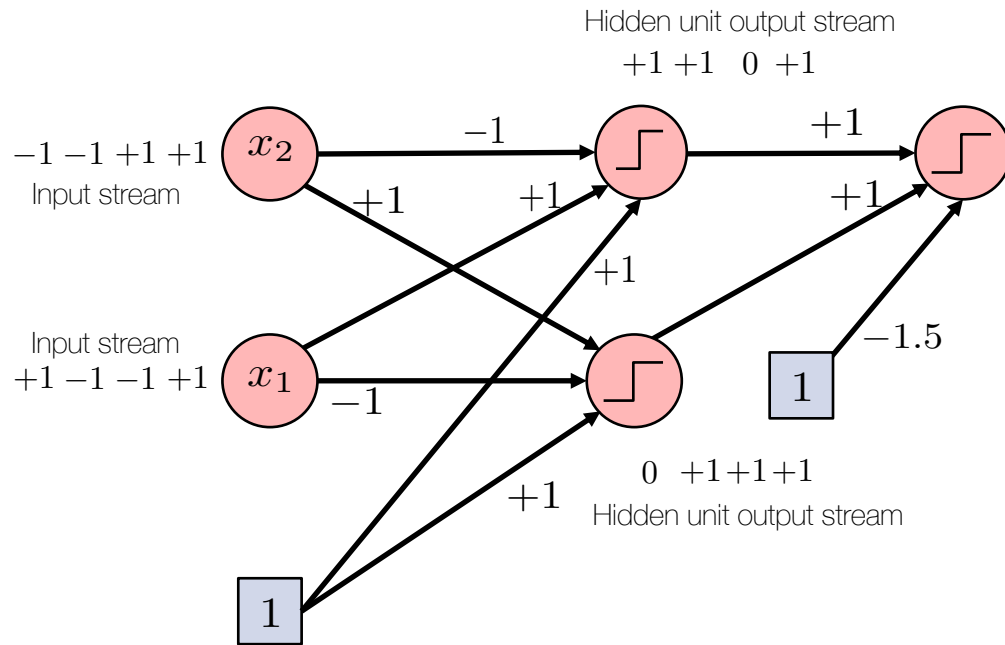
○ Cons

- We don't exactly know the nonlinear feature transform function
- If there are too many data samples the size of the kernel matrix is too large to handle
- Since we don't have the access to the mapping function, handling new data points is not straightforward
 - Though we can still do an out-of-samples extension
 - Basically the kernel matrix needs to grow



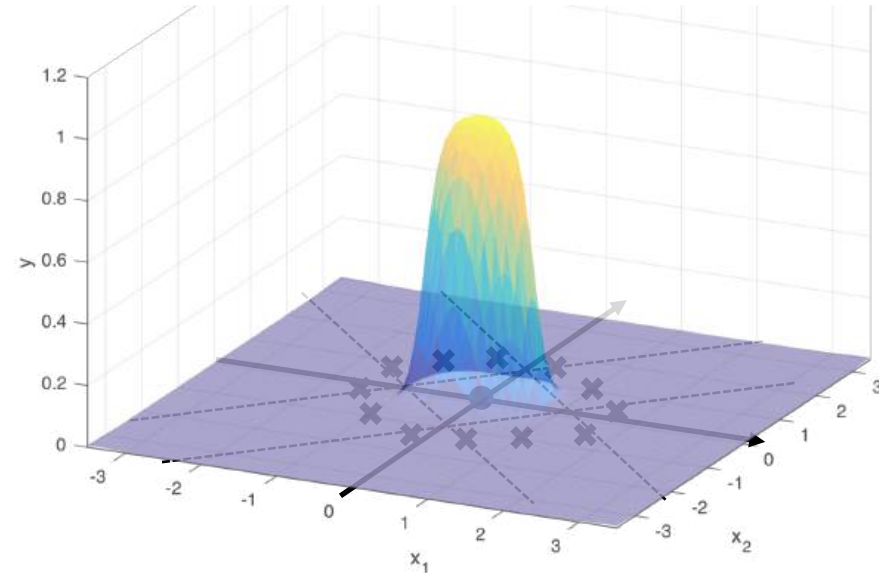
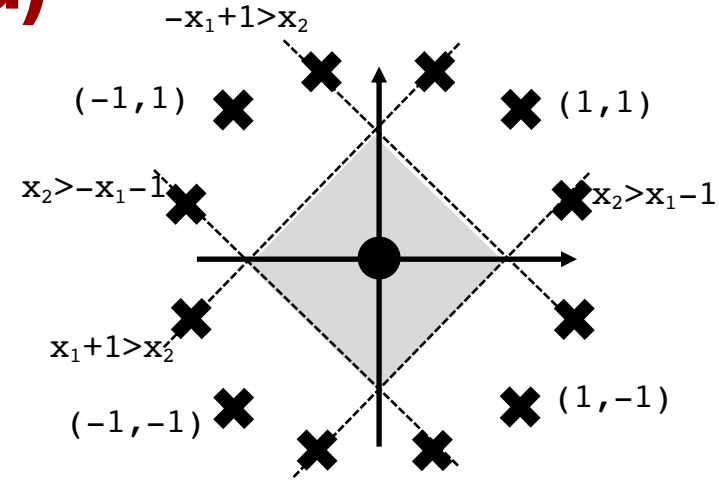
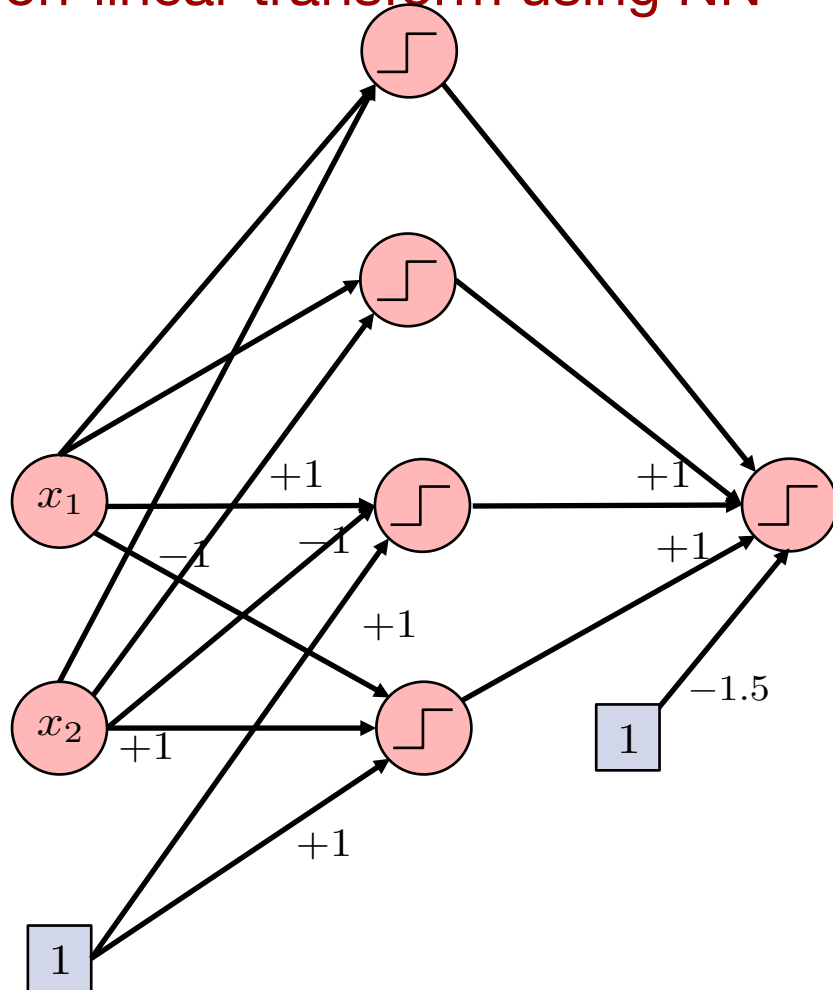
Neural Networks (Revisited)

- Non-linear transform using NN



Neural Networks (Revisited)

- Non-linear transform using NN



Multilayer Perceptron

- Backpropagation

- First we randomly initialize the weights and forwardpropagate

- Layer 1

$$\begin{bmatrix} z_2^{(1)} \\ z_1^{(1)} \end{bmatrix} = \begin{bmatrix} A_{22}^{(1)} & A_{21}^{(1)} & A_{20}^{(1)} \\ A_{12}^{(1)} & A_{11}^{(1)} & A_{10}^{(1)} \end{bmatrix} \cdot \begin{bmatrix} x_2^{(1)} \\ x_1^{(1)} \\ 1 \end{bmatrix} \quad \mathbf{z}^{(1)} = \mathbf{A}^{(1)} \mathbf{x}^{(1)}$$

$$\begin{bmatrix} x_2^{(2)} \\ x_1^{(2)} \end{bmatrix} = g \left(\begin{bmatrix} z_2^{(1)} \\ z_1^{(1)} \end{bmatrix} \right) \quad \mathbf{x}^{(2)} = g(\mathbf{z}^{(1)})$$

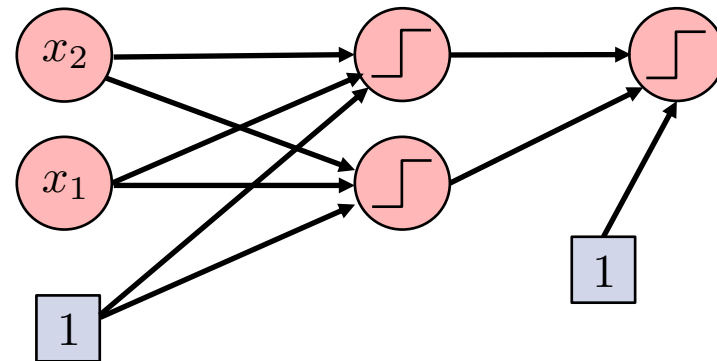
- Final layer

$$z_1^{(2)} = \begin{bmatrix} A_{12}^{(2)} & A_{11}^{(2)} & A_{10}^{(2)} \end{bmatrix} \cdot \begin{bmatrix} x_2^{(2)} \\ x_1^{(2)} \\ 1 \end{bmatrix} \quad z^{(2)} = \mathbf{A}^{(2)} \mathbf{x}^{(2)}$$

$$\hat{y} = g(z^{(2)})$$

- We can think of

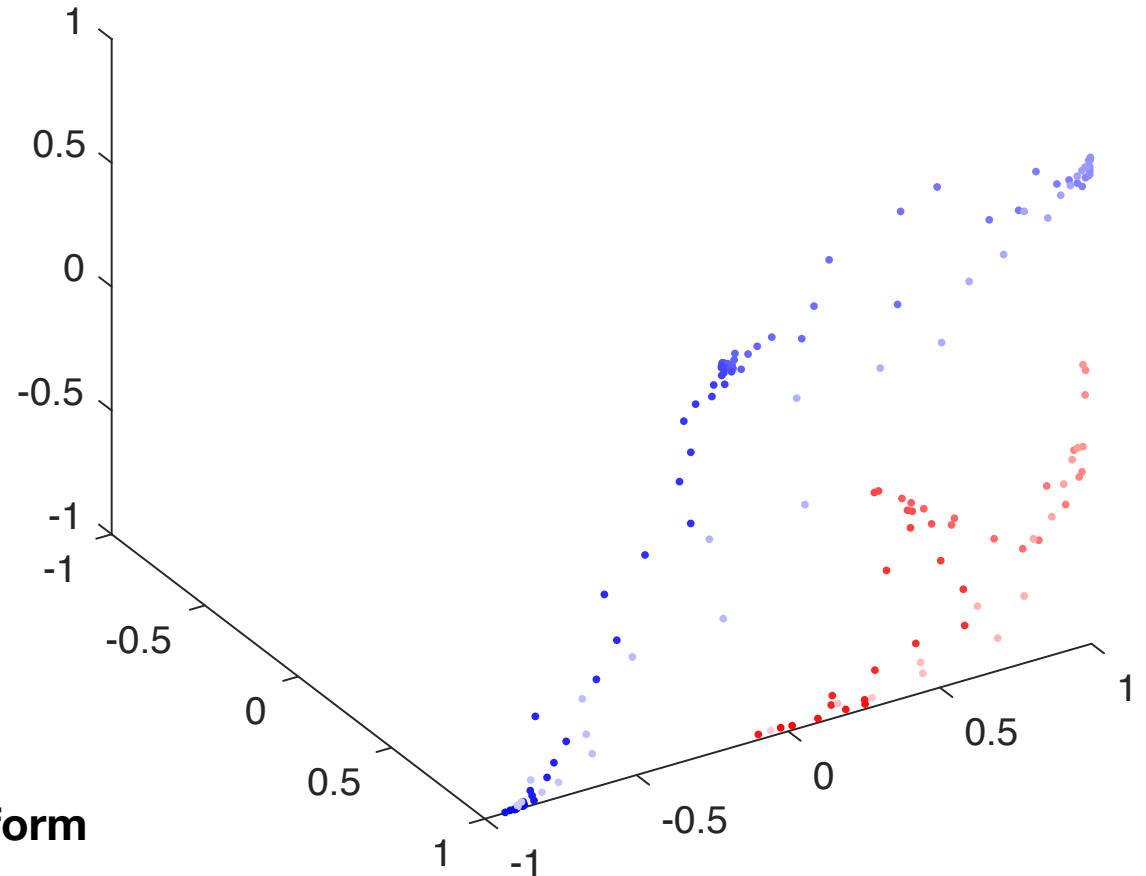
- The first layer as a nonlinear feature transform procedure
 - The final layer as the actual mapping function approximation layer (e.g. classification)
 - So what is the variable that represents the features?



Neural Networks (Revisited)

- Non-linear transform using NN

- The network takes the concentric circles as the input
- Plot of $x^{(2)}$
 - The output of the first layer
 - There are three hidden units
- They are now linearly separable
- **NN can do some nonlinear feature transform**



Reading

- Chapter 6.7, 15.2.4, 6.3





Thank You!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING