

Human Motion Analysis

Lecture 3: Dimensionality reduction

Raquel Urtasun

TTI Chicago

March 8, 2010

Contents of today's lecture

- How to deal with high-dimensional data.
- We will talk about different dimensionality reduction techniques
 - Linear models: PCA, CCA, etc.
 - Graph based methods: Isomap, Locally linear embedding, laplacian eigenmaps, etc.
 - Latent variable models: GTM and GPLVM
- We will see some examples in practice.

Materials used for this lecture

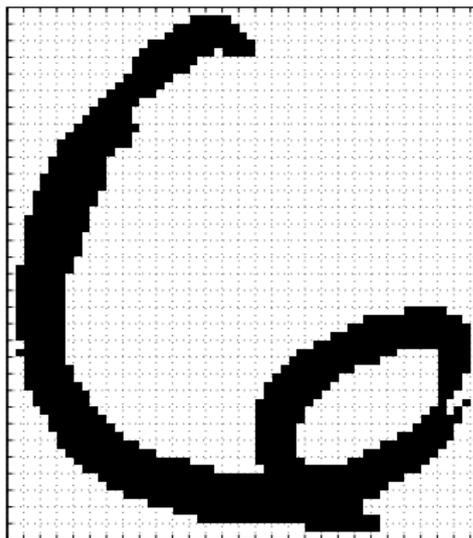
This lecture is based on two tutorials

- The ICML 2009 tutorial on dimensionality reduction given by Neil Lawrence.
- The tutorial on dimensionality reduction that Carl Ek gave at Oxford a few years back.

Thanks Neil and Carl for your slides!

USPS Data Set Handwritten Digit

- 3648 Dimensions
 - 64 rows by 57 columns
 - Space contains more than just this digit.



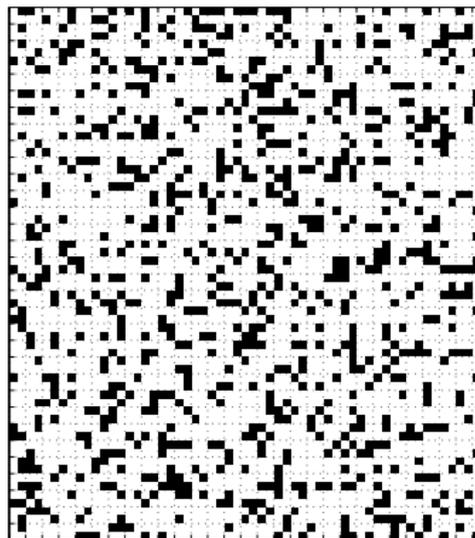
USPS Data Set Handwritten Digit

- 3648 Dimensions
 - 64 rows by 57 columns
 - Space contains more than just this digit.
 - Even if we sample every nanosecond from now until the end of the universe, you won't see the original six!



USPS Data Set Handwritten Digit

- 3648 Dimensions
 - 64 rows by 57 columns
 - Space contains more than just this digit.
 - Even if we sample every nanosecond from now until the end of the universe, you won't see the original six!

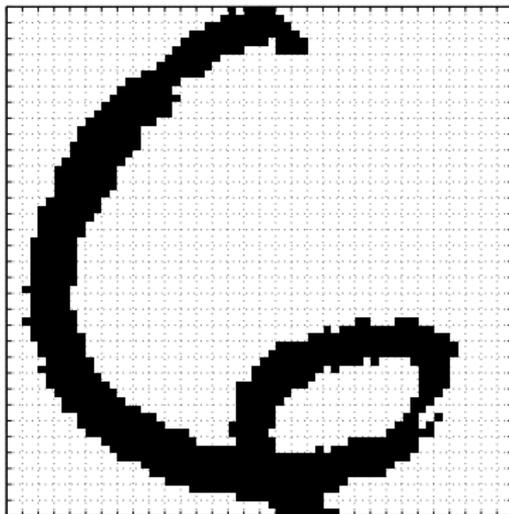


USPS Data Set Handwritten Digit

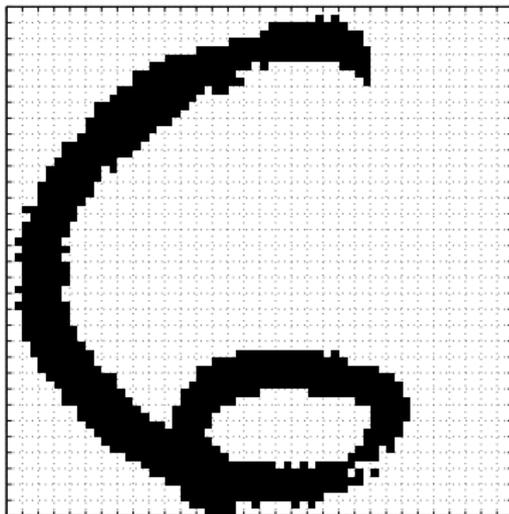
- 3648 Dimensions
 - 64 rows by 57 columns
 - Space contains more than just this digit.
 - Even if we sample every nanosecond from now until the end of the universe, you won't see the original six!



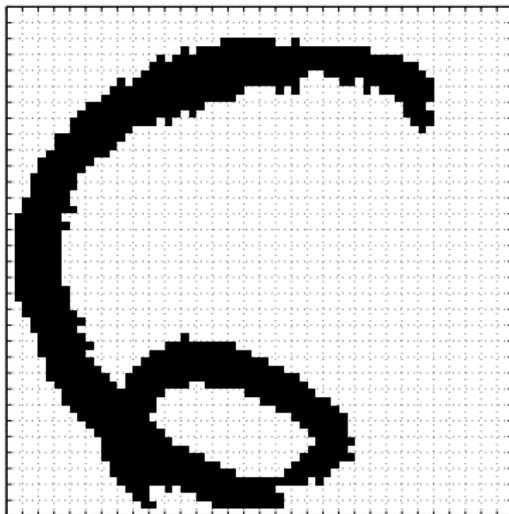
Rotate a 'Prototype'



Rotate a 'Prototype'

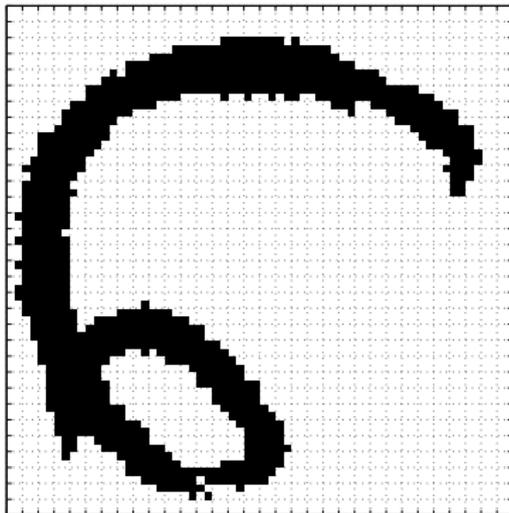


Rotate a 'Prototype'

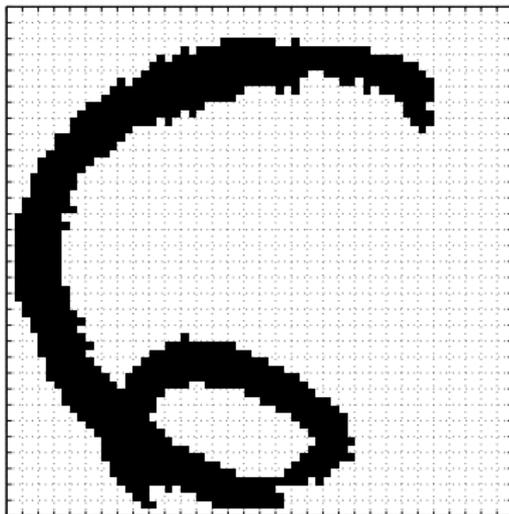


Simple model of a digit

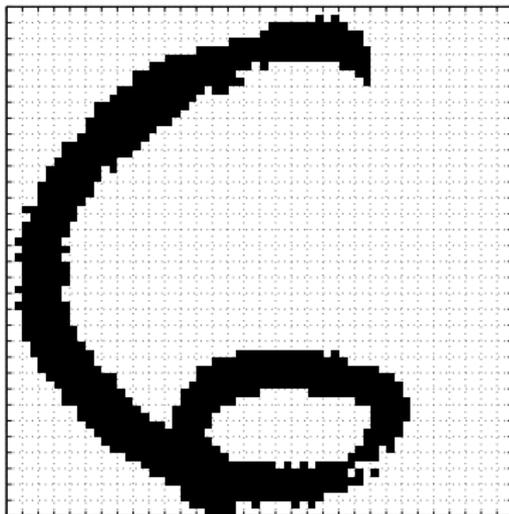
Rotate a 'Prototype'



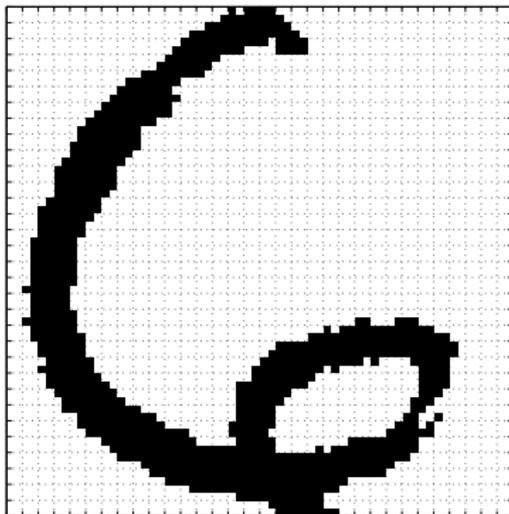
Rotate a 'Prototype'



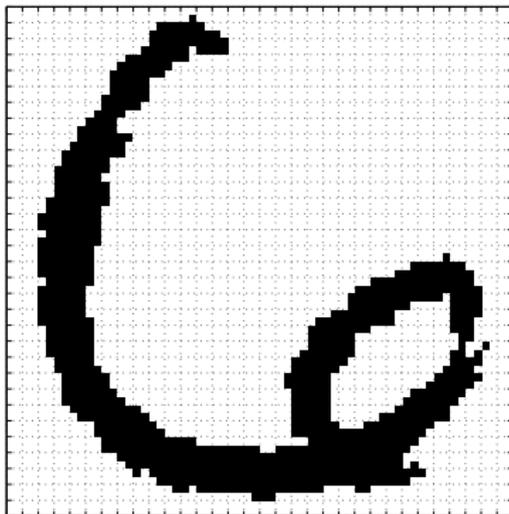
Rotate a 'Prototype'



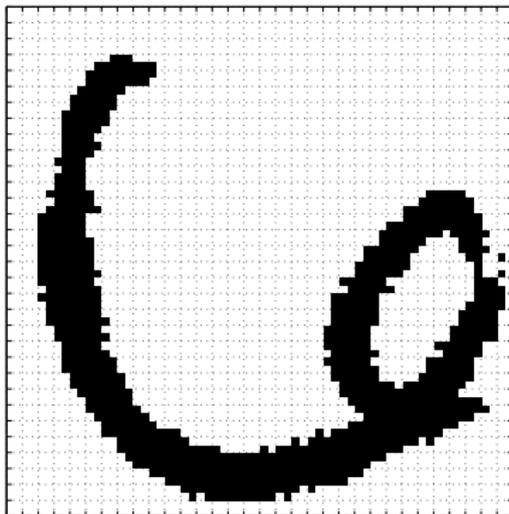
Rotate a 'Prototype'



Rotate a 'Prototype'

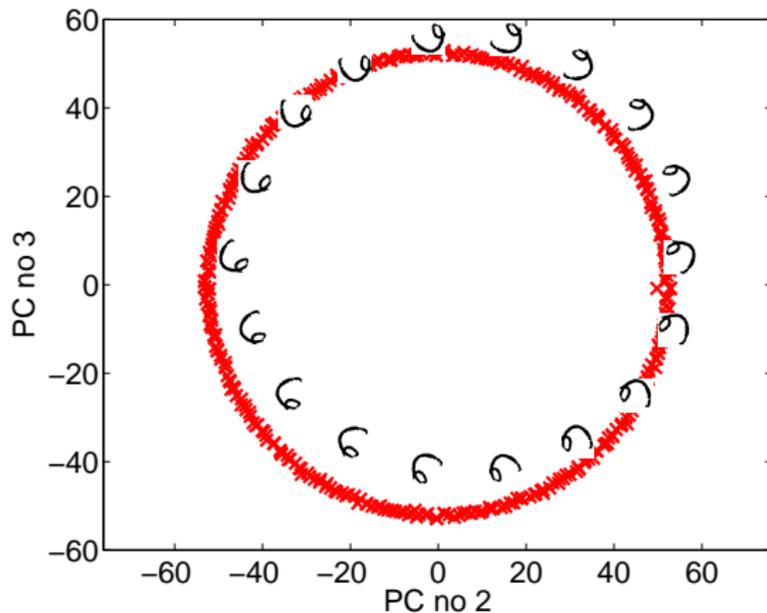


Rotate a 'Prototype'



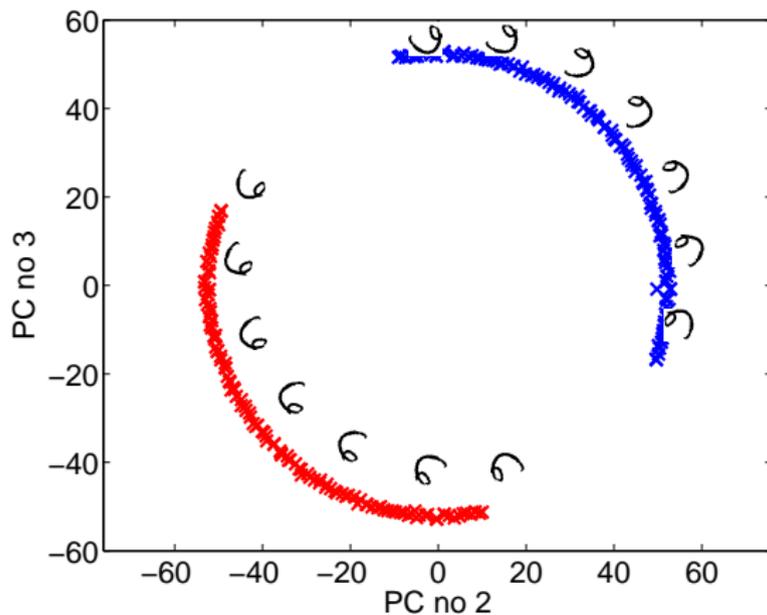
Two dimensional representation

```
demDigitsManifold[1 2], 'all')
```



Two dimensional representation

```
demDigitsManifold([1 2], 'sixnine')
```



Pure Rotation is too Simple

- In practice the data may undergo several distortions.
 - e.g. digits undergo 'thinning', translation and rotation.
- For data with 'structure':
 - we expect fewer distortions than dimensions;
 - we therefore expect the data to live on a lower dimensional manifold.
- Conclusion: deal with high dimensional data by looking for lower dimensional embedding.

q — dimension of latent/embedded space

D — dimension of data space

N — number of data points

centred data, $\mathbf{Y} = [\mathbf{y}_{1,:}, \dots, \mathbf{y}_{N,:}]^T = [\mathbf{y}_{:,1}, \dots, \mathbf{y}_{:,D}] \in \mathbb{R}^{N \times D}$
latent variables, $\mathbf{X} = [\mathbf{x}_{1,:}, \dots, \mathbf{x}_{N,:}]^T = [\mathbf{x}_{:,1}, \dots, \mathbf{x}_{:,q}] \in \mathbb{R}^{N \times q}$
mapping matrix, $\mathbf{W} \in \mathbb{R}^{D \times q}$

$\mathbf{a}_{i,:}$ is a vector from the i th row of a given matrix \mathbf{A}

$\mathbf{a}_{:,j}$ is a vector from the j th row of a given matrix \mathbf{A}

\mathbf{X} and \mathbf{Y} are *design matrices*

- Data covariance given by $N^{-1}\mathbf{Y}^T\mathbf{Y}$

$$\text{cov}(\mathbf{Y}) = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_{i,:} \mathbf{y}_{i,:}^T = \frac{1}{N} \mathbf{Y}^T \mathbf{Y}$$

- Inner product matrix given by $\mathbf{Y}\mathbf{Y}^T$

$$\mathbf{K} = (k_{i,j})_{i,j}, \quad k_{i,j} = \mathbf{y}_{i,:}^T \mathbf{y}_{j,:}$$

Types of approaches

- Linear dimensionality reduction
- Graph-based methods: based on preserving geodesic distances
- Non linear Latent variable models

Linear Dimensionality Reduction

- Two dimensional plane projected into a three dimensional space.

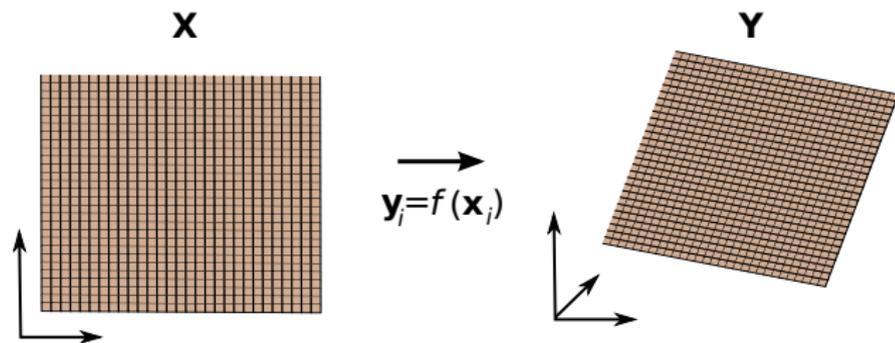


Figure: Mapping a 2D plane to a higher dimensional space in a linear way.

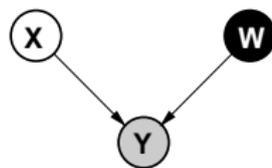
Linear Latent Variable Model

- Represent data, **Y**, with a lower dimensional set of latent variables **X**.
- Assume a linear relationship of the form

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\eta}_{i,:}, \quad \text{where } \boldsymbol{\eta}_{i,:} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

Probabilistic PCA

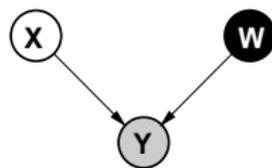
- Linear-Gaussian relationship between latent variables and data.
- \mathbf{X} are 'nuisance' variables.



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(y_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

Probabilistic PCA

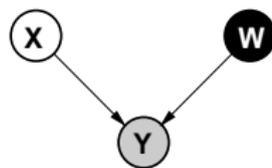
- Linear-Gaussian relationship between latent variables and data.
- \mathbf{X} are 'nuisance' variables.
- Latent variable model approach:



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(y_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

Probabilistic PCA

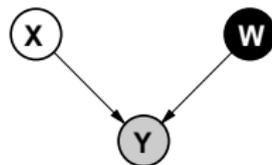
- Linear-Gaussian relationship between latent variables and data.
- \mathbf{X} are 'nuisance' variables.
- Latent variable model approach:
 - Define Gaussian prior over *latent space*, \mathbf{X} .



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(y_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

Probabilistic PCA

- Linear-Gaussian relationship between latent variables and data.
- \mathbf{X} are 'nuisance' variables.
- Latent variable model approach:
 - Define Gaussian prior over *latent space*, \mathbf{X} .
 - Integrate out nuisance *latent variables*.



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(y_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_{i,:} | \mathbf{0}, \mathbf{I})$$

Probabilistic PCA

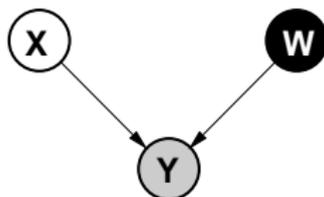
- Linear-Gaussian relationship between latent variables and data.
- \mathbf{X} are 'nuisance' variables.
- Latent variable model approach:
 - Define Gaussian prior over *latent space*, \mathbf{X} .
 - Integrate out nuisance *latent variables*.

$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_{i,:} | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

Probabilistic PCA Max. Likelihood Soln (Tipping and Bishop, 1999b)



$$p(\mathbf{Y}|\mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$$

Probabilistic PCA Max. Likelihood Soln (Tipping and Bishop, 1999b)

$$p(\mathbf{Y}|\mathbf{W}) = \prod_{j=1}^D \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{C}), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$$

$$\log p(\mathbf{Y}|\mathbf{W}) = -\frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \text{tr}(\mathbf{C}^{-1}\mathbf{Y}^T\mathbf{Y}) + \text{const.}$$

If \mathbf{U}_q are first q principal eigenvectors of $N^{-1}\mathbf{Y}^T\mathbf{Y}$ and the corresponding eigenvalues are Λ_q ,

$$\mathbf{W} = \mathbf{U}_q\mathbf{L}\mathbf{R}^T, \quad \mathbf{L} = (\Lambda_q - \sigma^2\mathbf{I})^{\frac{1}{2}}$$

where \mathbf{R} is an arbitrary rotation matrix.

- Very similar to PCA, but with a more complex notion of noise:

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \epsilon$$

with $E\{\epsilon\epsilon^T\} = \Sigma$.

- If the noise is known, then the factors can be estimated using PCA of a modified matrix

$$\mathbf{C} - \Sigma$$

with \mathbf{C} the covariance matrix of the data.

- If the noise is not known, then there exist different algorithms in the literature to solve this.
- We will not see them in this class.

Why Probabilistic PCA?

- What is the point in probabilistic methods?
- Could we not just project with regular PCA?

Why Probabilistic PCA?

- What is the point in probabilistic methods?
- Could we not just project with regular PCA?
 - Integration within other models (e.g. mixtures of PCA (Tipping and Bishop, 1999a) , temporal models).

Why Probabilistic PCA?

- What is the point in probabilistic methods?
- Could we not just project with regular PCA?
 - Integration within other models (e.g. mixtures of PCA (Tipping and Bishop, 1999a) , temporal models).
 - Model selection through Bayesian treatment of parameters (Bishop 1999) .

Why Probabilistic PCA?

- What is the point in probabilistic methods?
- Could we not just project with regular PCA?
 - Integration within other models (e.g. mixtures of PCA (Tipping and Bishop, 1999a) , temporal models).
 - Model selection through Bayesian treatment of parameters (Bishop 1999) .
 - Marginalisation of missing data (Tipping and Bishop, 1999b) .

Why Probabilistic PCA?

- What is the point in probabilistic methods?
- Could we not just project with regular PCA?
 - Integration within other models (e.g. mixtures of PCA (Tipping and Bishop, 1999a) , temporal models).
 - Model selection through Bayesian treatment of parameters (Bishop 1999) .
 - Marginalisation of missing data (Tipping and Bishop, 1999b) .

Why Probabilistic PCA?

- What is the point in probabilistic methods?
- Could we not just project with regular PCA?
 - Integration within other models (e.g. mixtures of PCA (Tipping and Bishop, 1999a) , temporal models).
 - Model selection through Bayesian treatment of parameters (Bishop 1999) .
 - Marginalisation of missing data (Tipping and Bishop, 1999b) .

Note: These same advantages hold for Factor Analysis

- Distributions can behave very non-intuitively in high dimensions.
- Fortunately, most data is not really high dimensional.
- Probabilistic PCA exploits linear low dimensional structure in the data.
 - Probabilistic interpretation brings with it many advantages: extensibility, Bayesian approaches, missing data.
- We will now motivate the need for *non linear* dimensionality reduction.

Why non-linear dimensionality reduction?

- Complex datasets cannot be represented linearly.

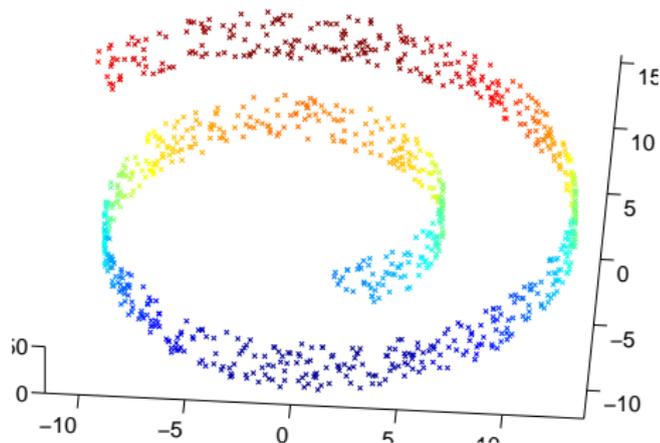


Figure: The 'Swiss Roll' data set is data in three dimensions that is inherently two dimensional.

- We will see non-linear latent variable models and spectral methods.

Spectral Approaches

- Classical Multidimensional Scaling (MDS) (Mardia et al. 1979) .
 - Uses eigenvectors of similarity matrix.
- Kernel PCA (Scholkopf et al., 1998)
 - Provides a representation and a mapping — representation is high dimensional though!
 - Mapping is implied through the use of a kernel function as a similarity matrix.
- Isomap (Tenenbaum et al., 2000) is MDS with a particular proximity measure.
 - Approximate distances measures along the manifold.
 - Compute neighborhood and compute shortest distance in graph.
 - Use classical MDS on that distance matrix.

Non Probabilistic Existing Methods II

- **Locally Linear Embedding** (Roweis and Saul, 2000) .
 - Looks to preserve locally linear relationships in a low dimensional space.
 - Compute neighborhood and point find reduced dimensional relationships that preserve local linearity.
- **Laplacian Eigenmaps** (Belkin and Niyogi, 2003) .
 - Uses spectral graph theory and information geometric arguments to form embedding.
 - Compute neighborhood, graph Laplacian and seek 2nd lowest eigenvector.
- **Maximum Variance Unfolding** (Weinberger et al., 2004) .
 - Compute neighborhood, constrain local distances to be preserved.
 - Maximise the variance in latent space.

Iterative Methods

- Multidimensional Scaling (MDS)
 - Iterative optimisation of a stress function (Kruskal, 1964) .
- Sammon Mappings (Sammon, 1969) .
 - Strictly speaking not a mapping — similar to iterative MDS.
- NeuroScale (Lowe and Tipping, 1997)
 - Augmentation of iterative MDS methods with a mapping.

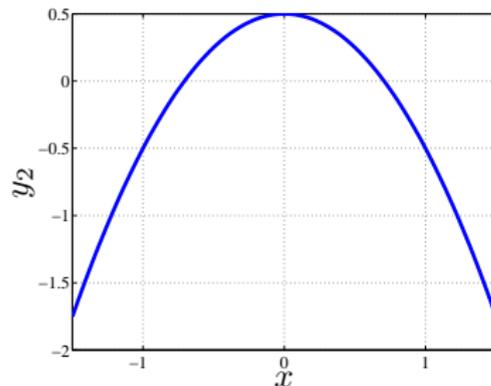
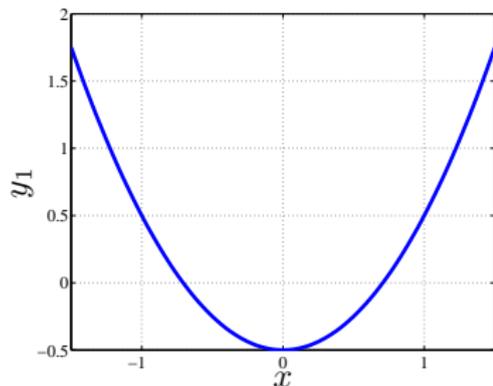
Local Distance Preservation

- Most of the above dimensional reduction techniques preserve local distances.
 - Probabilistic Approaches do not.
- Probabilistic approaches map smoothly from latent to data space.
 - Points close in latent space are close in data space.
 - This does not imply points close in data space are close in latent space.
- Spectral approaches map smoothly from data to latent space.
 - Points close in data space are close in latent space.
 - This does not imply points close in latent space are close in data space.

Forward Mapping

- Mapping from 1-D latent space to 2-D data space.

$$y_1 = x^2 - 0.5, \quad y_2 = -x^2 + 0.5$$

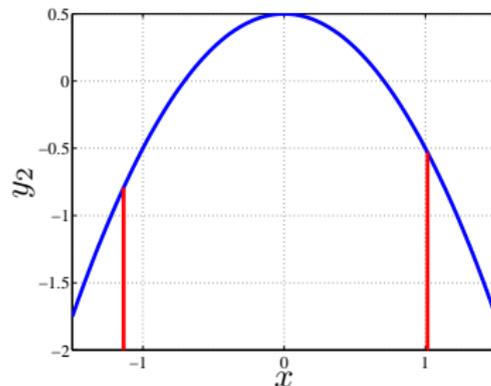
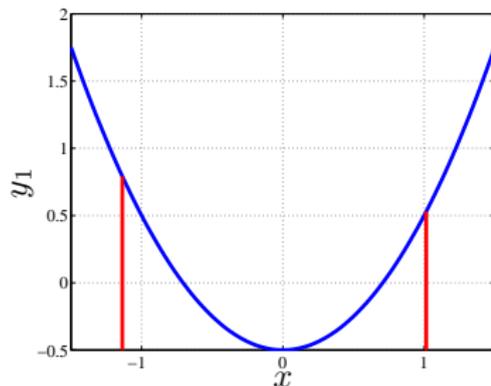


Distance Preservation

Forward Mapping

- Mapping from 1-D latent space to 2-D data space.

$$y_1 = x^2 - 0.5, \quad y_2 = -x^2 + 0.5$$

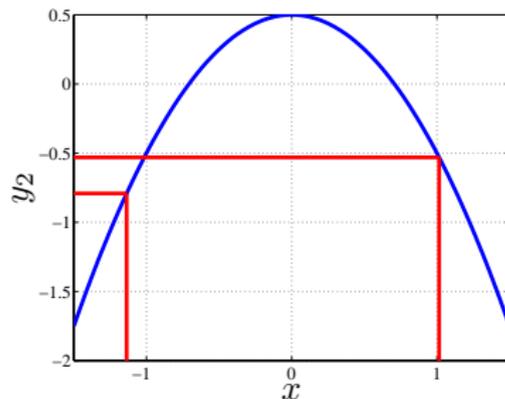
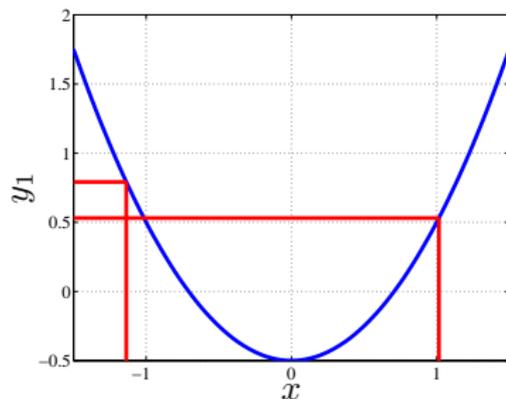


Distance Preservation

Forward Mapping

- Mapping from 1-D latent space to 2-D data space.

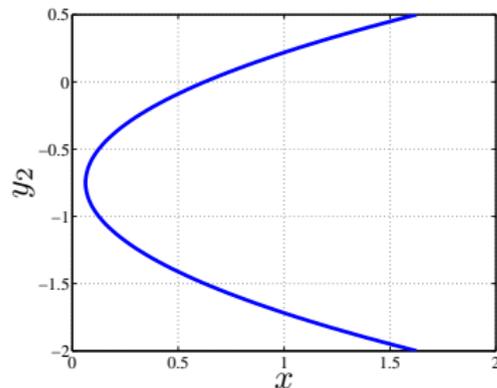
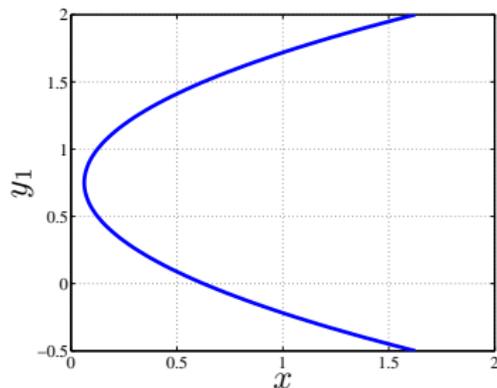
$$y_1 = x^2 - 0.5, \quad y_2 = -x^2 + 0.5$$



Backward Mapping

- Mapping from 2-D data space to 1-D latent.

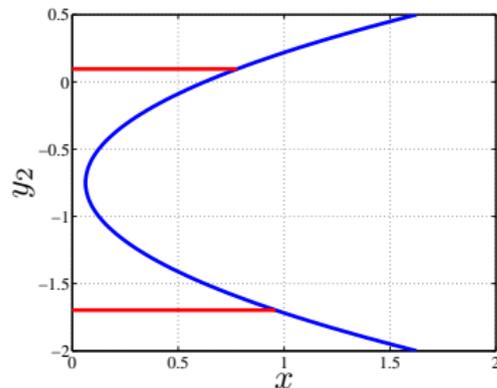
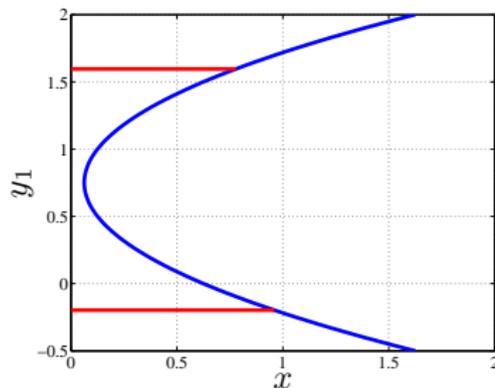
$$x = 0.5 (y_1^2 + y_2^2 + 1)$$



Backward Mapping

- Mapping from 2-D data space to 1-D latent.

$$x = 0.5 (y_1^2 + y_2^2 + 1)$$

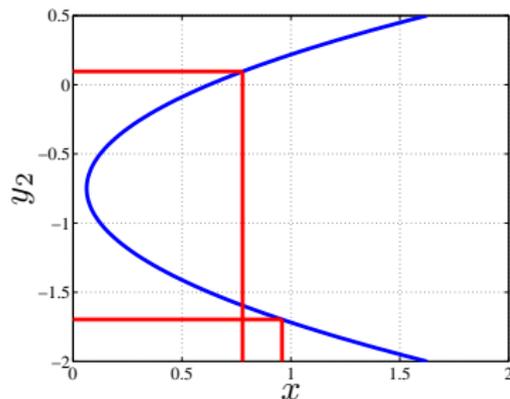
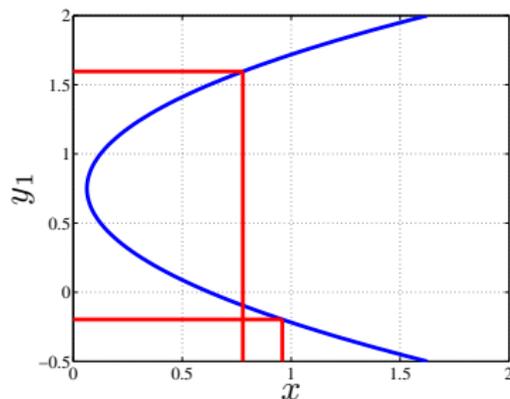


Distance Preservation

Backward Mapping

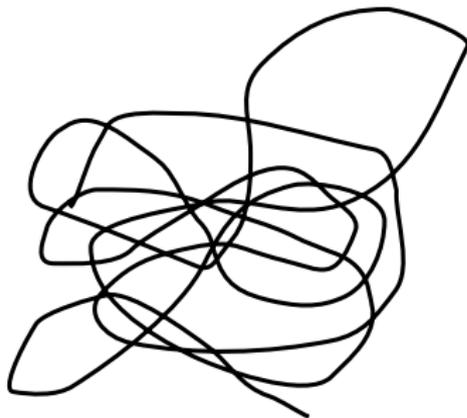
- Mapping from 2-D data space to 1-D latent.

$$x = 0.5 (y_1^2 + y_2^2 + 1)$$



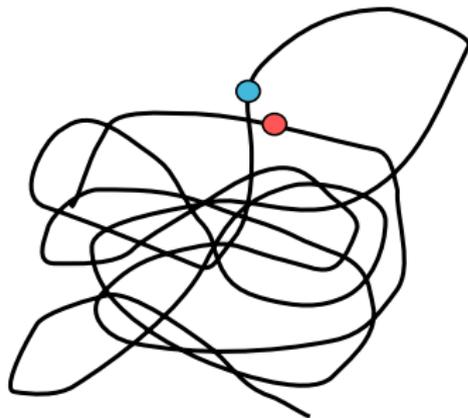
Tangled String

- Sometimes local distance preservation in data space is wrong.
- The pink and blue ball should be separated.



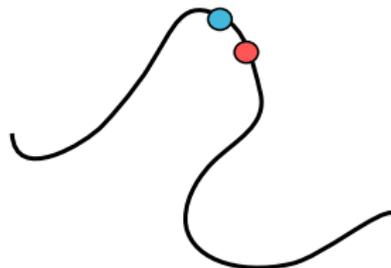
Tangled String

- Sometimes local distance preservation in data space is wrong.
- The pink and blue ball should be separated.
- But the assumption makes the problem simpler (for spectral methods it is convex).



Tangled String

- Sometimes local distance preservation in data space is wrong.
- The pink and blue ball should be separated.
- But the assumption makes the problem simpler (for spectral methods it is convex).



Good

- Unique optimum.

But

- Non trivial for dealing with missing data.
- Difficult to extend (e.g. temporal data) in a principled way.

We are going to see in more detail:

- Multidimensional Scaling (MDS)
- Kernel PCA
- Isomap
- Maximum Variance Unfolding (MVU)
- Locally Linear Embedding (LLE)
- Laplacian Eigenmaps

- Classical statistical approach: represent via proximities (Mardia, 1972).
- Proximity data: similarities or dissimilarities.
- Example of a dissimilarity matrix: a *distance matrix*.

$$d_{i,j} = \|\mathbf{y}_{i,:} - \mathbf{y}_{j,:}\|_2 = \sqrt{(\mathbf{y}_{i,:} - \mathbf{y}_{j,:})^\top (\mathbf{y}_{i,:} - \mathbf{y}_{j,:})}$$

- For a data set can display as a matrix.

Interpoint Distances for Rotated Sixes

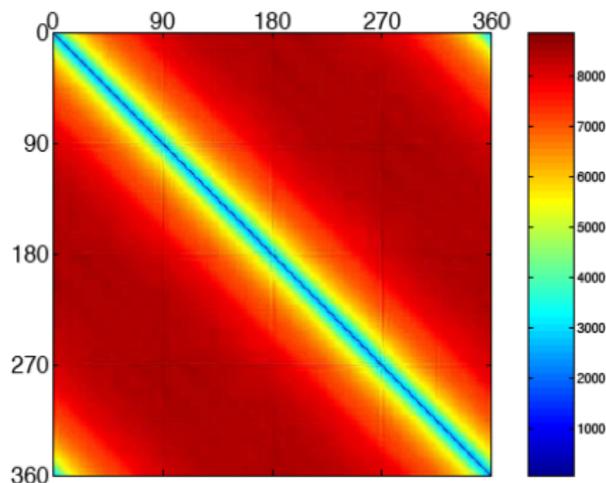


Figure: Interpoint distances for the rotated digits data.

Multidimensional Scaling

- Find a configuration of points, \mathbf{X} , such that each

$$\delta_{i,j} = \|\mathbf{x}_{i,:} - \mathbf{x}_{j,:}\|_2$$

closely matches the corresponding $d_{i,j}$ in the distance matrix.

- Need an objective function for matching $\mathbf{\Delta} = (\delta_{i,j})_{i,j}$ to $\mathbf{D} = (d_{i,j})_{i,j}$.

- An entrywise L_1 norm on difference between squared distances

$$E(\mathbf{X}) = \sum_{i=1}^N \sum_{j=1}^N |d_{ij}^2 - \delta_{ij}^2|.$$

- Reduce dimension by selecting features from data set.
- Select for \mathbf{X} , in turn, the column from \mathbf{Y} that most reduces this error until we have the desired q .
- To minimise $E(\mathbf{Y})$ we compose \mathbf{X} by extracting the columns of \mathbf{Y} which have the largest variance.

Reconstruction from Latent Space

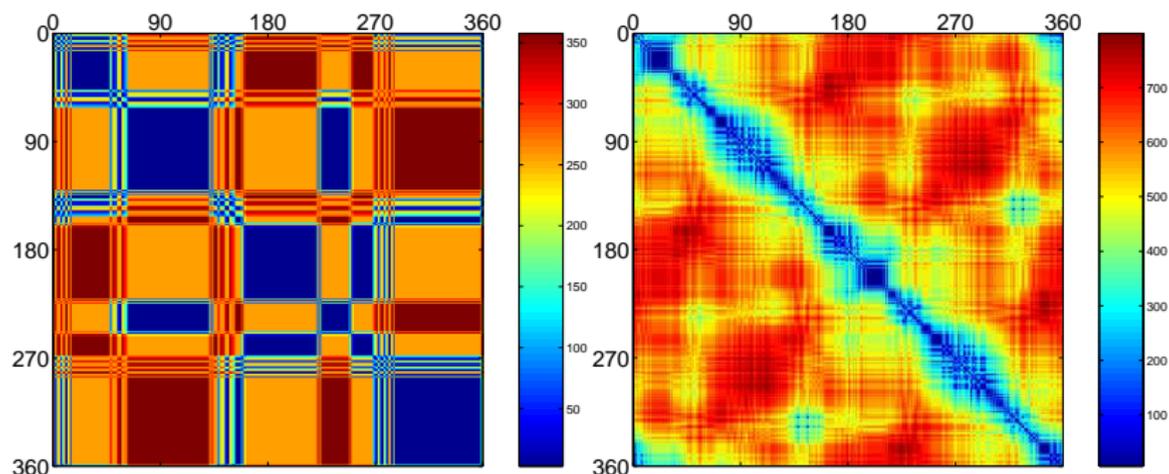


Figure:

Left: distances reconstructed with two dimensions. *Right:* distances reconstructed with 10 dimensions.

Reconstruction from Latent Space

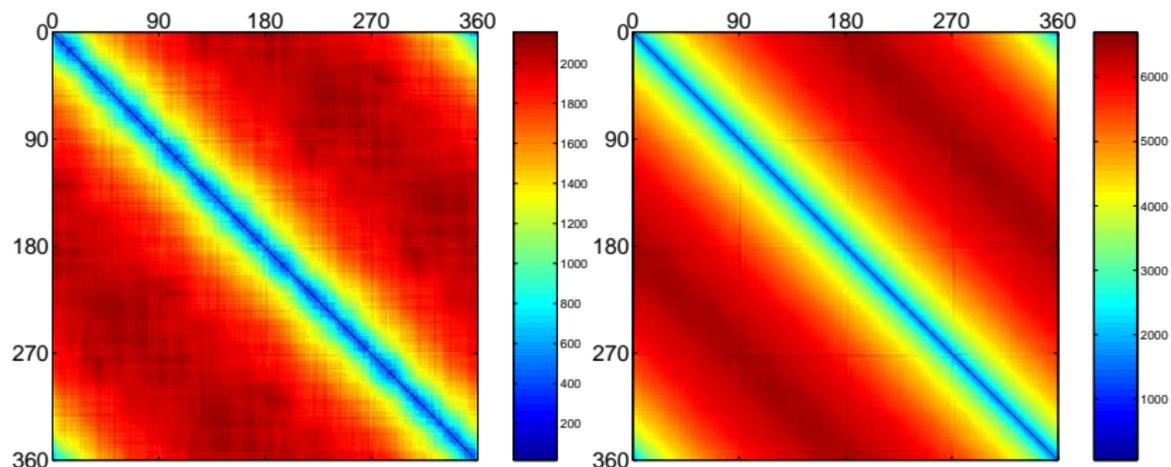


Figure:

Left: distances reconstructed with 100 dimensions. *Right:* distances reconstructed with 1000 dimensions.

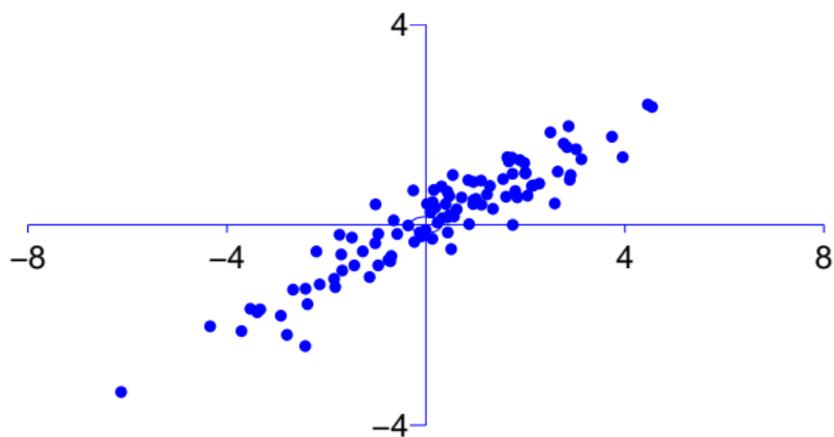


Figure: demRotationDist. Feature selection via distance preservation.

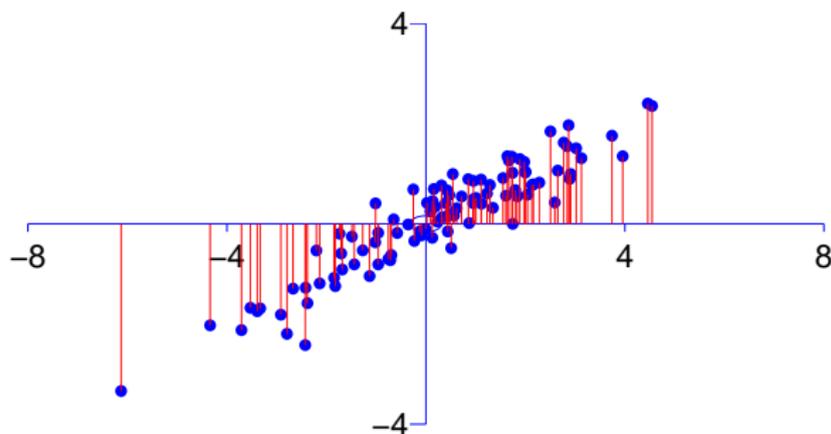


Figure: `demRotationDist`. Feature selection via distance preservation.

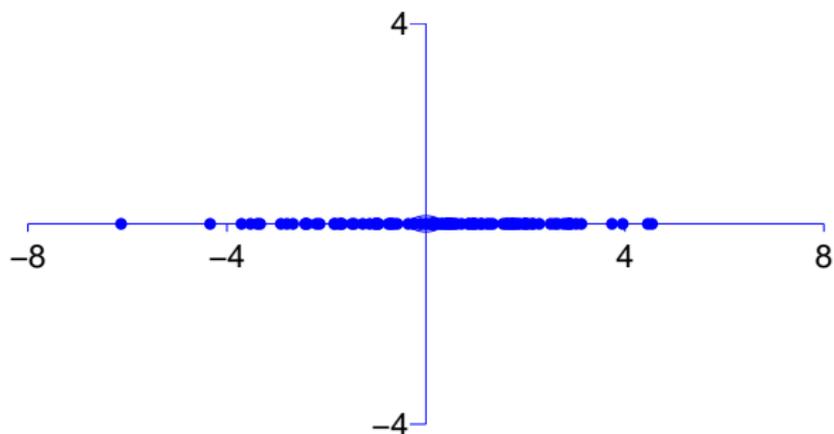


Figure: `demRotationDist`. Feature selection via distance preservation.

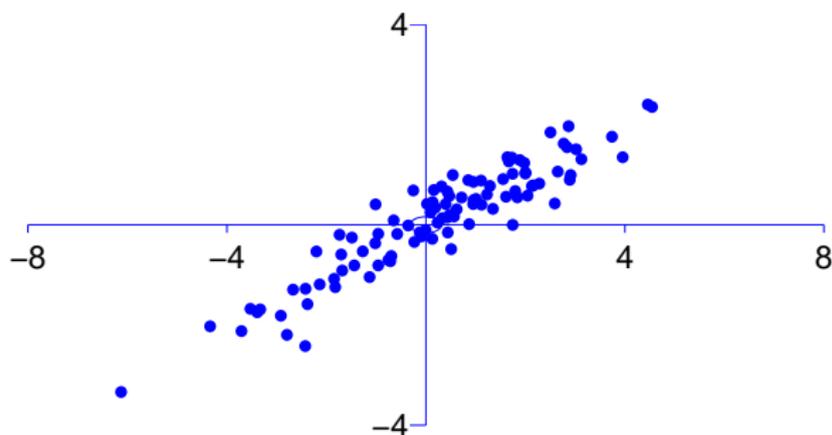


Figure: `demRotationDist`. Rotation preserves interpoint distances. .

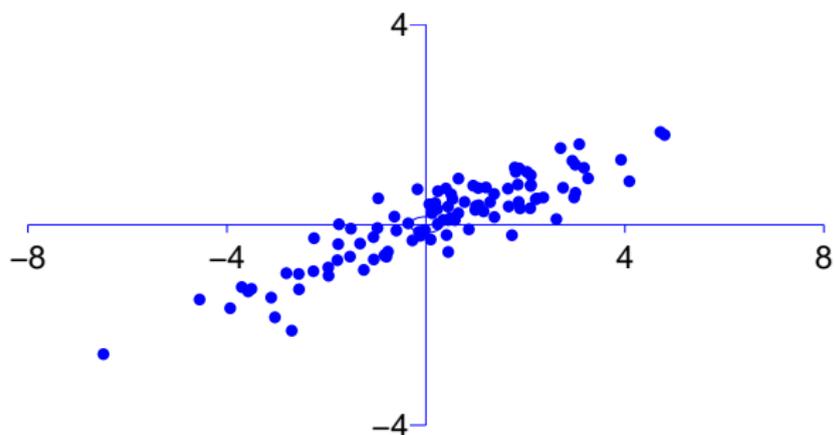


Figure: `demRotationDist`. Rotation preserves interpoint distances. .

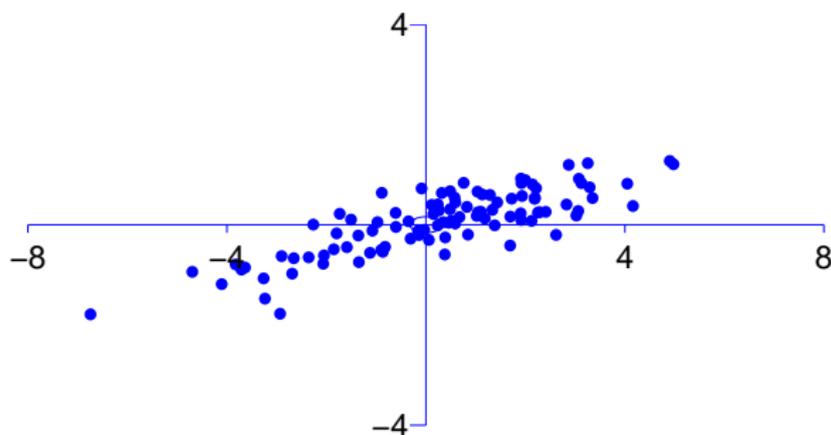


Figure: `demRotationDist`. Rotation preserves interpoint distances. .

Feature Extraction

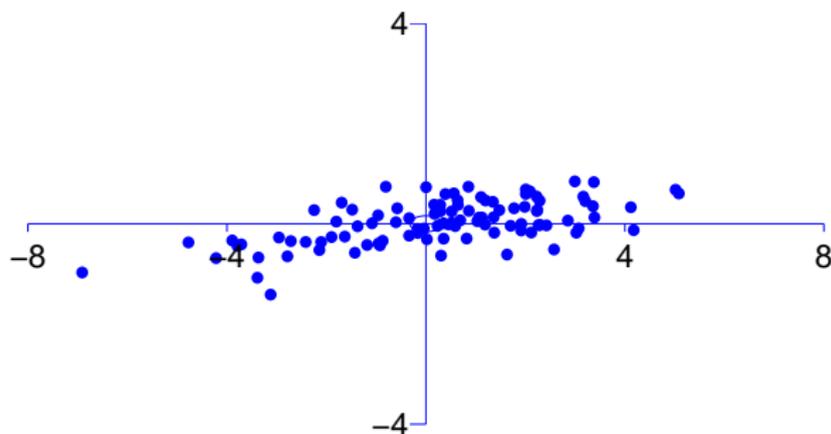


Figure: `demRotationDist`. Rotation preserves interpoint distances. .

Feature Extraction

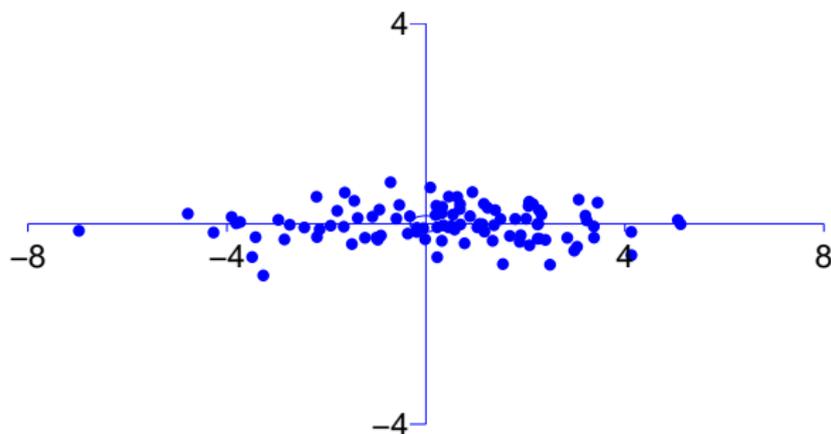


Figure: `demRotationDist`. Rotation preserves interpoint distances. .

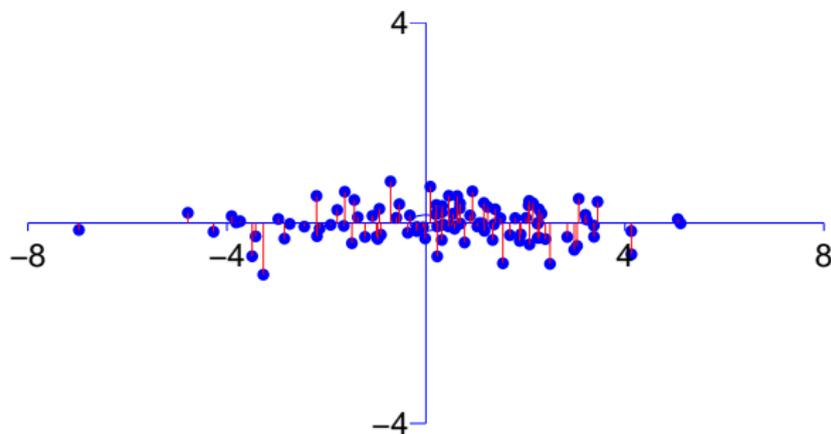


Figure: `demRotationDist`. Rotation preserves interpoint distances. Residuals are much reduced.

Feature Extraction

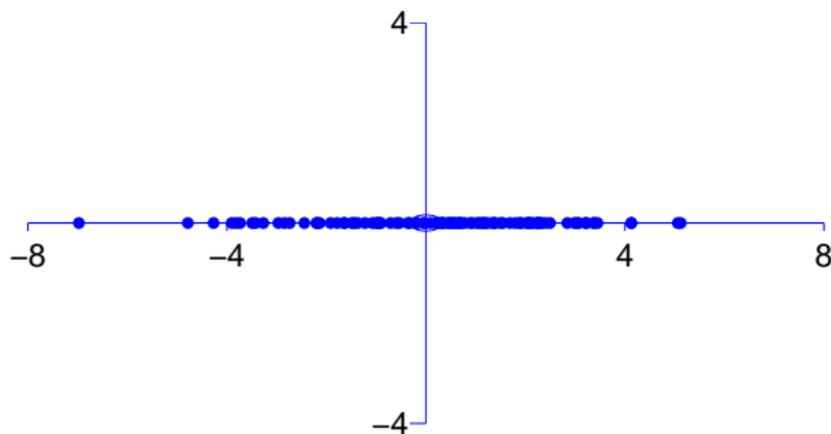


Figure: `demRotationDist`. Rotation preserves interpoint distances. Residuals are much reduced.

Which Rotation?

- We need the rotation that will minimise residual error.
- We already derived an algorithm for discarding directions.
- Discard direction with *maximum variance*.
- Error is then given by the sum of residual variances.

$$E(\mathbf{X}) = 2N^2 \sum_{k=q+1}^D \sigma_k^2.$$

- Rotations of data matrix *do not* effect this analysis.

Rotation Reconstruction from Latent Space

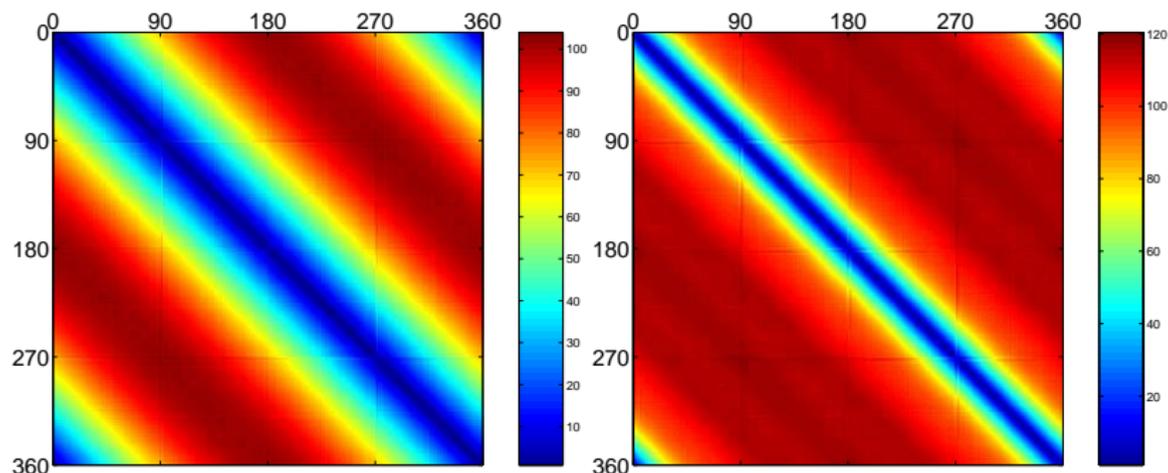


Figure:

Left: distances reconstructed with two dimensions. *Right:* distances reconstructed with 10 dimensions.

Rotation Reconstruction from Latent Space

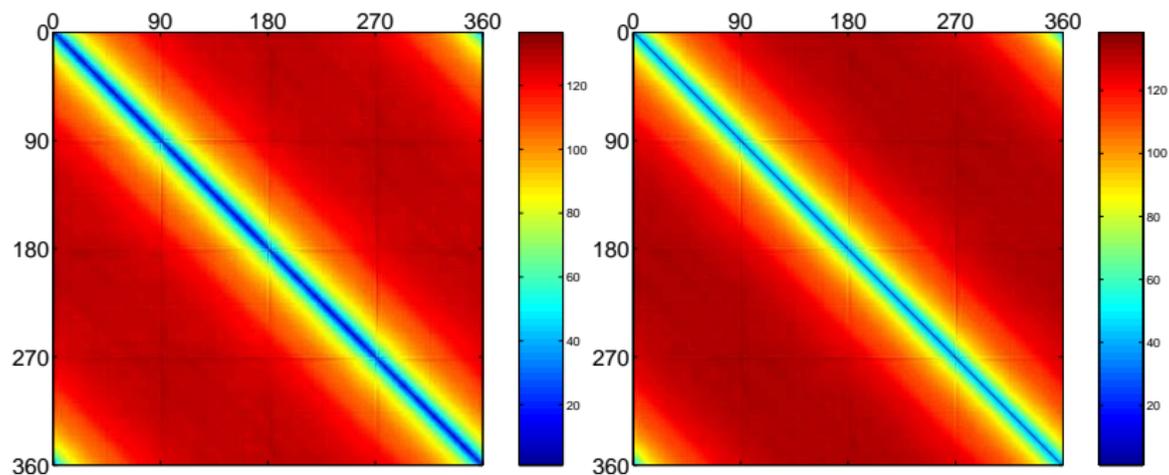


Figure:

Left: distances reconstructed with 100 dimensions. *Right:* distances reconstructed with 360 dimensions.

Reminder: Principal Component Analysis

- How do we find these directions?
- Find directions in data with maximal variance.
 - That's what PCA does!
- **PCA**: rotate data to extract these directions.
- **PCA**: work on the sample covariance matrix $\mathbf{S} = N^{-1}\hat{\mathbf{Y}}^T\hat{\mathbf{Y}}$.

Distance to Similarity: Gaussian Covariances

- Translate between covariance and distance.
 - Consider a vector sampled from a zero mean Gaussian distribution,

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}).$$

- Expected square distance between two elements of this vector is

$$d_{i,j}^2 = \langle (z_i - z_j)^2 \rangle$$

$$d_{i,j}^2 = \langle z_i^2 \rangle + \langle z_j^2 \rangle - 2 \langle z_i z_j \rangle$$

under a zero mean Gaussian with covariance given by \mathbf{K} this is

$$d_{i,j}^2 = k_{i,i} + k_{j,j} - 2k_{i,j}.$$

Take the distance to be square root of this,

$$d_{i,j} = (k_{i,i} + k_{j,j} - 2k_{i,j})^{\frac{1}{2}}.$$

Standard Transformation

- This transformation is known as the *standard transformation* between a similarity and a distance (Mardia et al. pg 402, 1979) .
- If the covariance is of the form $\mathbf{K} = \hat{\mathbf{Y}}\hat{\mathbf{Y}}^\top$ then $k_{i,j} = \mathbf{y}_{i,:}^\top \mathbf{y}_{j,:}$ and

$$d_{i,j} = \left(\mathbf{y}_{i,:}^\top \mathbf{y}_{i,:} + \mathbf{y}_{j,:}^\top \mathbf{y}_{j,:} - 2\mathbf{y}_{i,:}^\top \mathbf{y}_{j,:} \right)^{\frac{1}{2}} = \|\mathbf{y}_{i,:} - \mathbf{y}_{j,:}\|_2 .$$

- For other distance matrices this gives us an approach to convert to a similarity matrix or kernel matrix so we can perform classical MDS.

Example: Road Distances with Classical MDS

- Classical example: redraw a map from road distances (see e.g. Mardia et al. 1979).
- Here we use distances across Europe.
 - Between each city we have road distance.
 - Enter these in a distance matrix.
 - Convert to a similarity matrix using the covariance interpretation.
 - Perform eigendecomposition.

Other Distance Similarity Measures

- Can use similarity/distance of your choice.
- Beware though!
 - The similarity must be positive semi definite for the distance to be Euclidean.
 - Why? Can immediately see positive definite is sufficient from the “covariance interpretation” .
 - For more details see (Mardia et al. 1979, Theorem 14.2.2) .

- All Mercer kernels are positive semi definite.
- Example, squared exponential (also known as RBF or Gaussian)

$$k_{i,j} = \exp \left(-\frac{\|\mathbf{y}_{i,:} - \mathbf{y}_{j,:}\|^2}{2l^2} \right).$$

This leads to a kernel eigenvalue problem.

- This is known as Kernel PCA Scholkopf et al. 1998.

Implied Distances on Rotated Sixes

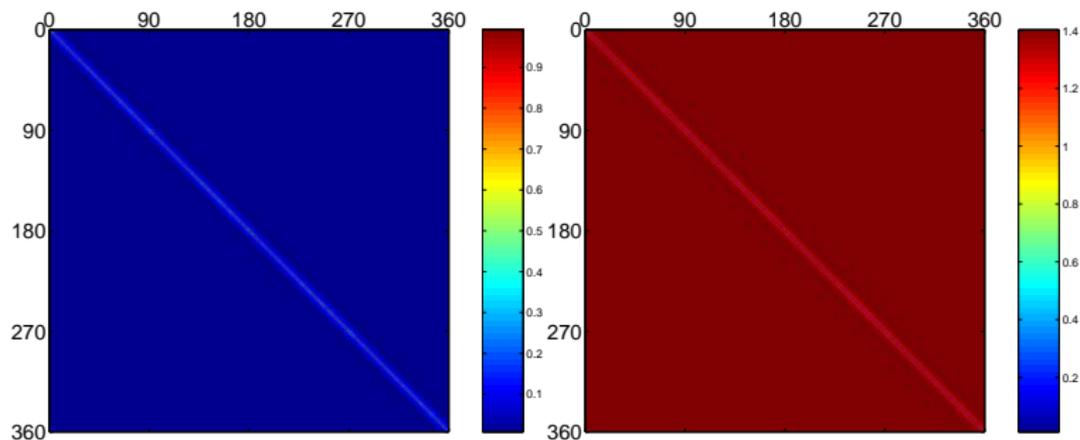


Figure: *Left:* similarity matrix for RBF kernel on rotated sixes. *Right:* implied distance matrix for kernel on rotated sixes. Note that most of the distances are set to $\sqrt{2} \approx 1.41$.

Kernel PCA on Rotated Sixes

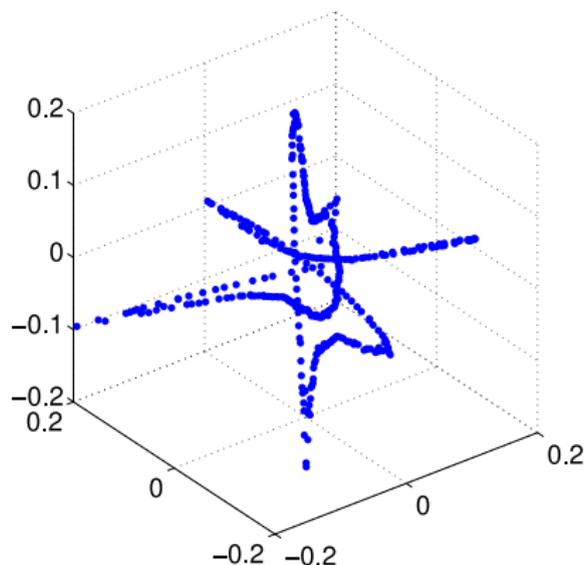


Figure: `demSixKpca`. The fifth, sixth and seventh dimensions of the latent space for kernel PCA. Points spread out along axes so that dissimilar points are always $\sqrt{2}$ apart.

- Multidimensional scaling: preserve a distance matrix.
- Classical MDS
 - a particular objective function
 - for Classical MDS distance matching is equivalent to maximum variance
 - spectral decomposition of the similarity matrix
- For Euclidean distances in \mathbf{Y} space classical MDS is equivalent to PCA.
 - known as principal coordinate analysis (PCO)
- Haven't discussed choice of distance matrix.

Non-Linear vs. Linear — Local vs. Global

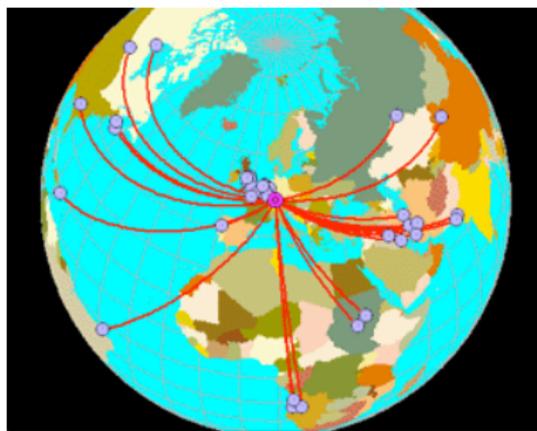
- MDS and PCA re-parametrise data based on **global** structures (linear) in the given representation of the data
- Idea: **Local** structure of given representation is close to the manifold structure
- Want to “unravel” local structure of data globally

Proximity Graph

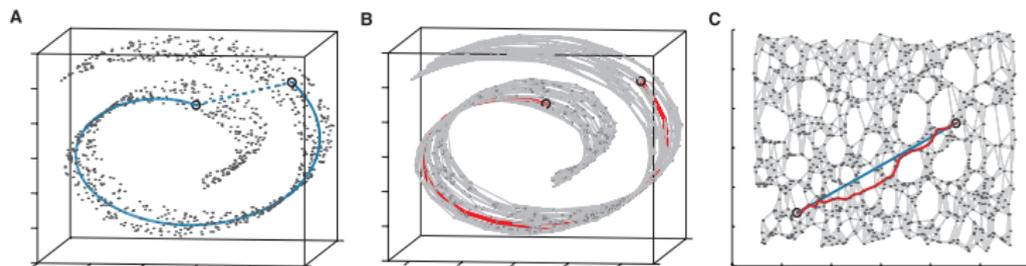
- 1 Identify neighbors of each data point $\mathbf{y}_i \in N(\mathbf{y}_j)$
- 2 Build graph $\mathbf{P} = \left\{ \underbrace{\mathbf{Y}}_{\text{vertexset}}, \underbrace{\mathbf{W}}_{\text{edgeset}} \right\}$
 - Put edges between vertices's in neighborhood
 - Assume \mathbf{P} connected (and in most cases symmetric)
- 3 **Objective:** Complete \mathbf{P} to make it fully connected
- 4 Different algorithms have different strategies
 - What are the edge weights?
 - How to complete \mathbf{P}

- *Tenenbaum, de Silva, Langford* - Science December 2000
- Local Proximity Graph
- Edge Weights Euclidean distances

- MDS finds geometric configuration preserving distances
- MDS applied to Manifold distance
- Geodesic Distance = Manifold Distance
- “Chicken and Egg” Cannot compute geodesic distance without knowing manifold



- Geodesic Distance can be approximated by shortest path through local proximity matrix
- Compute distance matrix by completing Proximity Graph



- 1 Compute Neighbor relations

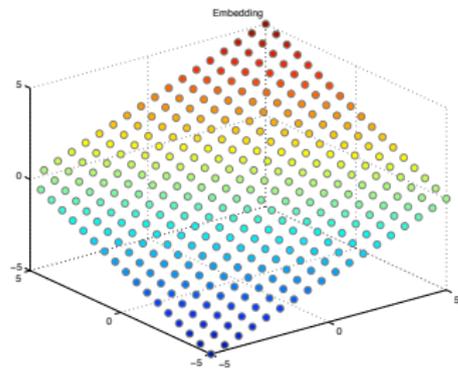
$$\Delta_{ij} = \begin{cases} \|\mathbf{y}_i - \mathbf{y}_j\|_2 & (\mathbf{y}_i, \mathbf{y}_j) \in W \\ \infty & \text{otherwise} \end{cases}$$

- 2 Complete Δ by Shortest path

$$\Delta_{ij} = \begin{cases} W_{ij} & (\mathbf{y}_i, \mathbf{y}_j) \in W \\ \text{shortestpath}(\mathbf{y}_i, \mathbf{y}_j, W) & \text{otherwise} \end{cases}$$

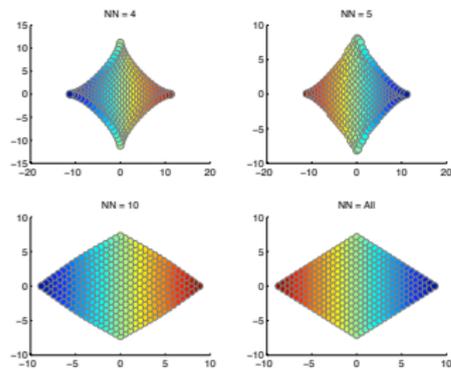
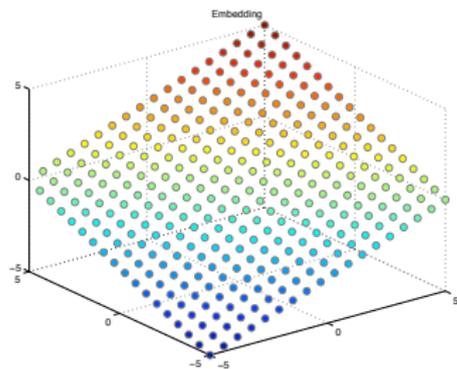
- 3 Apply MDS to Δ

Isomap: Example¹



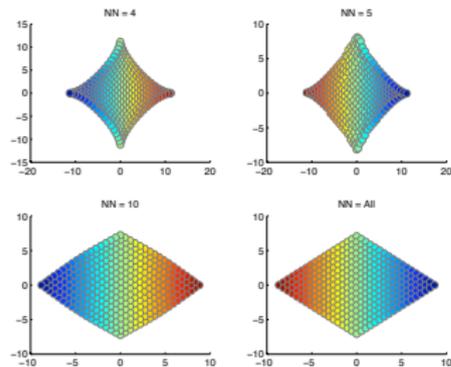
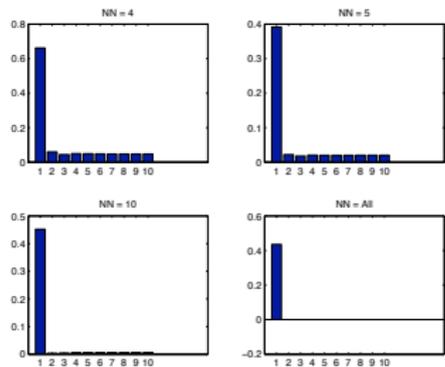
¹/algorithms/isomap_embed.m

Isomap: Example¹



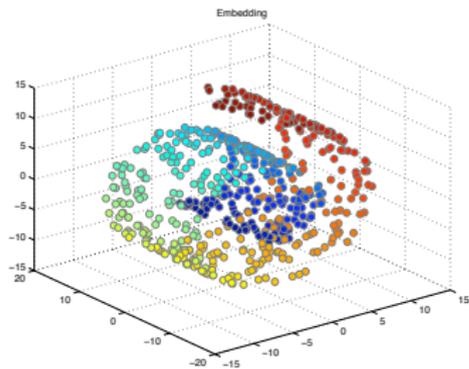
¹/algorithms/isomap_embed.m

Isomap: Example¹



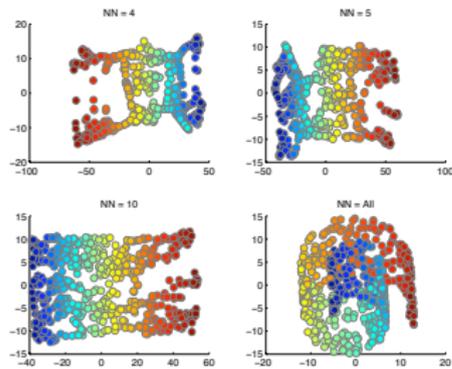
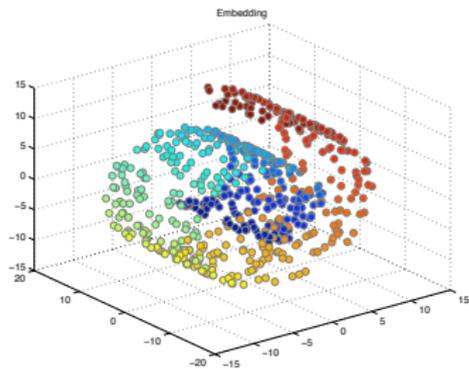
¹/algorithms/isomap_embed.m

Isomap: Example¹



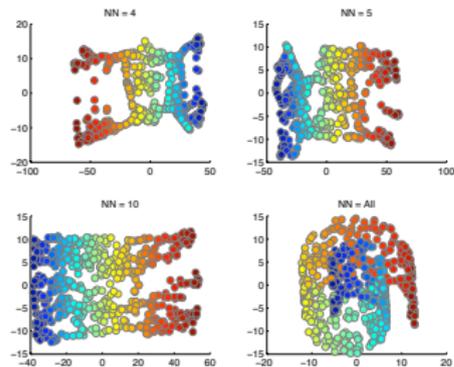
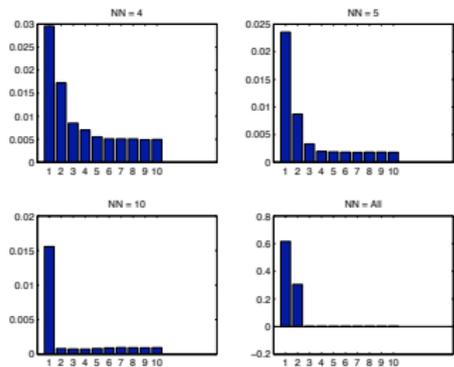
¹/algorithms/isomap_embed.m

Isomap: Example¹



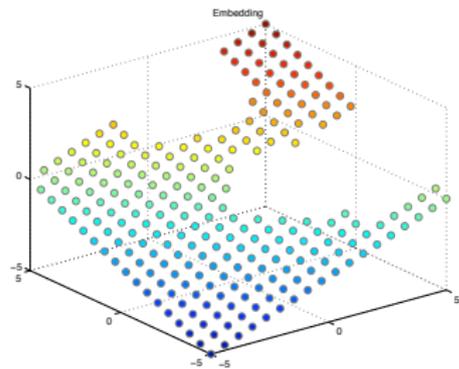
¹/algorithms/isomap_embed.m

Isomap: Example¹



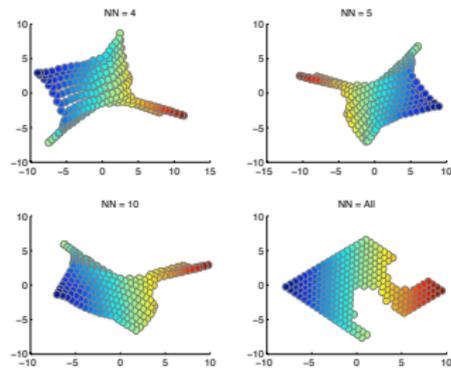
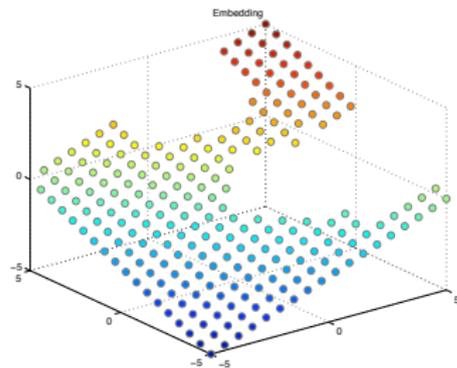
¹/algorithms/isomap_embed.m

Isomap: Example¹



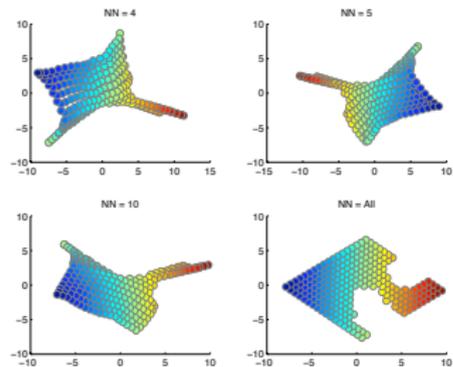
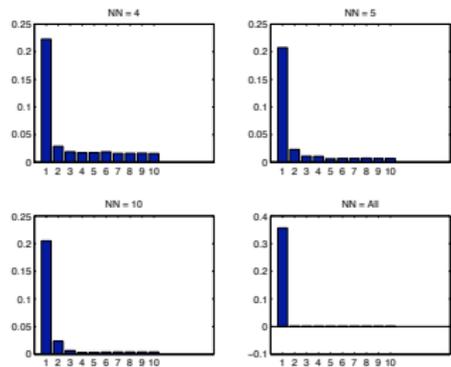
¹/algorithms/isomap_embed.m

Isomap: Example¹



¹/algorithms/isomap_embed.m

Isomap: Example¹



¹/algorithms/isomap_embed.m

- MDS on shortest path approximation of manifold distance
- + Simple
- + Intrinsic dimension from eigen spectra
 - Solves a very large eigenvalue problem
 - Cannot handle holes or non-convex manifold
 - Sensitive to “short circuit”
 - Increases rank of Gram matrix

Maximum Variance Unfolding

- *Weinberg, Sha, Saul* - ICML & CVPR 2004
- First presented as Semi-Definite Embeddings
- Formulate dimensionality reduction in terms of Gram matrix

Maximum Variance Unfolding

- Want to keep local structure $(\mathbf{y}_i, \mathbf{y}_j) \in W$

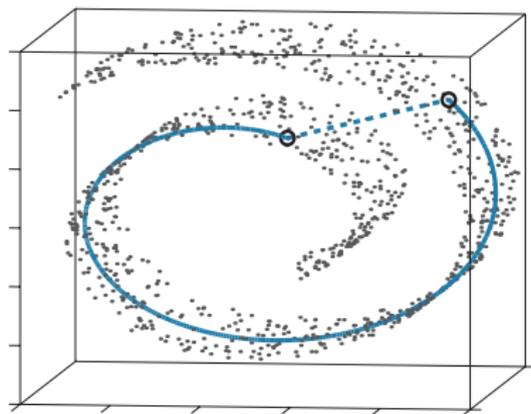
$$\begin{aligned} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 &= \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \\ \Rightarrow \mathbf{K}_{ii} + \mathbf{K}_{jj} - \mathbf{K}_{ij} - \mathbf{K}_{ji} &= \mathbf{G}_{ii} + \mathbf{G}_{jj} - \mathbf{G}_{ij} - \mathbf{G}_{ji} \end{aligned}$$

- Remove Translational Invariance

$$\left\| \sum_{i=1}^N \mathbf{x}_i \right\|_2^2 = 0 \Rightarrow \sum_{i=1}^N \sum_{j=1}^N \mathbf{K}_{ij} = 0$$

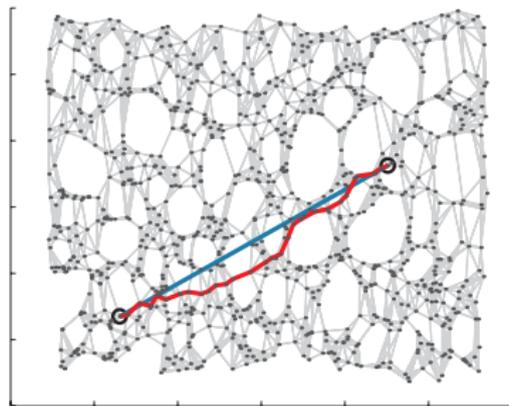
- Need to be valid Gram matrix $\Rightarrow \mathbf{K} \succcurlyeq 0$

Maximum Variance Unfolding



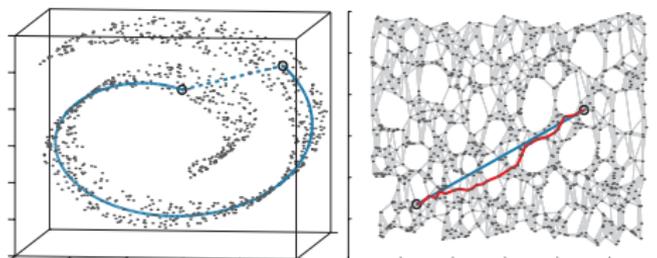
Any “fold” of the manifold between two points will **decrease** the *Euclidean* distance between the points while the *Manifold* distance remains **constant**

Maximum Variance Unfolding



If manifold is **maximally** stretched between two points the *Euclidean* distance will **equal** the *Manifold* distance

Maximum Variance Unfolding



Maximise all pairwise distance outside local neighborhood (upper bound)

$$\max \sum_{i=1}^N \sum_{j=1}^N \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$$
$$\Rightarrow \max(\text{trace}(\mathbf{K}))$$

Maximum Variance Unfolding: Algorithm

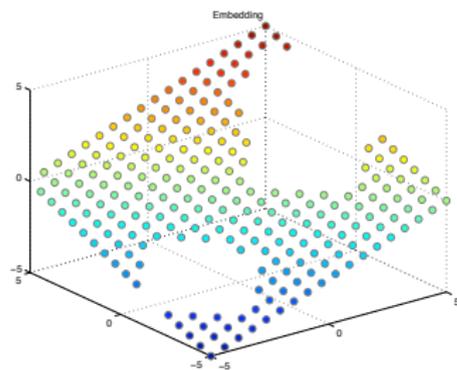
- 1 Compute Proximity Graph
- 2 Compute Local Gram Matrix \mathbf{G}
- 3 Compute Global Gram Matrix \mathbf{K}

$$\begin{aligned} & \max(\text{trace}(\mathbf{K})) \\ \text{subject to : } & \mathbf{K} \succcurlyeq 0 \\ & \sum_{i=1}^N \sum_{j=1}^N \mathbf{K}_{ij} = 0 \\ & \mathbf{K}_{ii} + \mathbf{K}_{jj} - \mathbf{K}_{ij} - \mathbf{K}_{ji} = \mathbf{G}_{ii} + \mathbf{G}_{jj} - \mathbf{G}_{ij} - \mathbf{G}_{ji} \end{aligned}$$

Instance of *Semidefinite Programming*

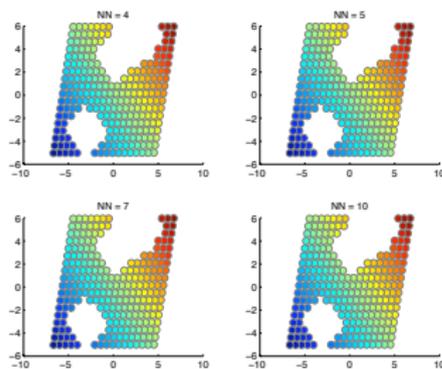
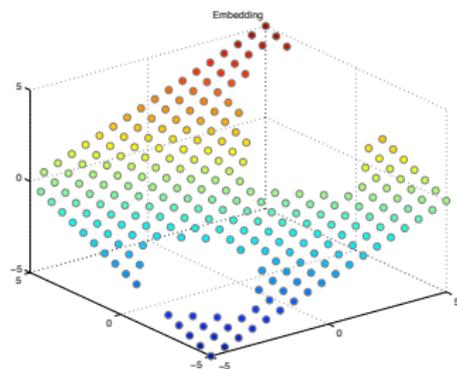
- 4 Apply MDS to \mathbf{K}

Maximum Variance Unfolding: Example²



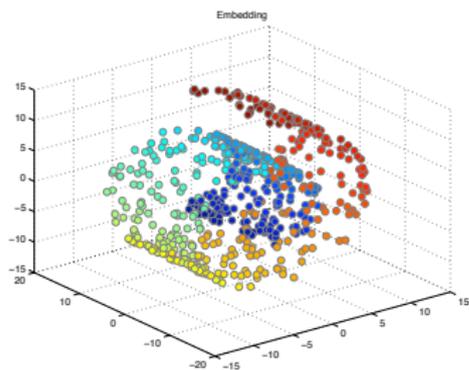
`2/algos/mvu_embed.m`

Maximum Variance Unfolding: Example²



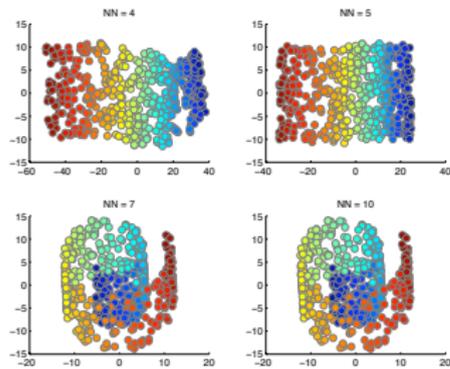
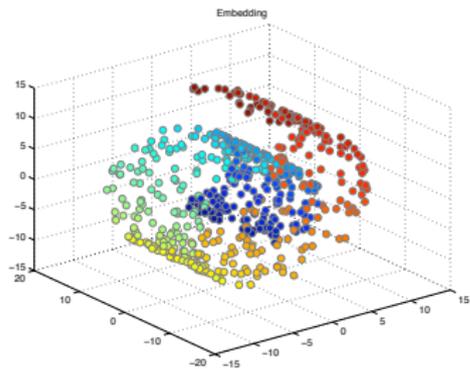
²/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example²



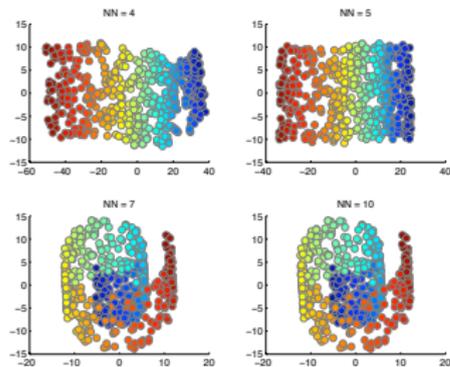
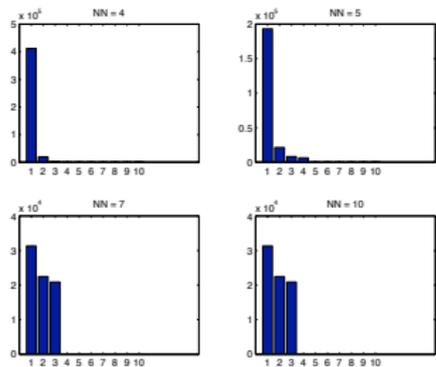
²/algos/mvu_embed.m

Maximum Variance Unfolding: Example²



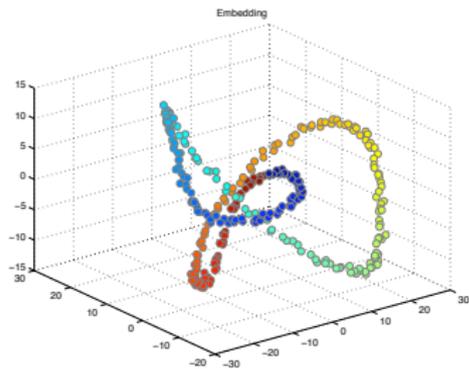
²/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example²



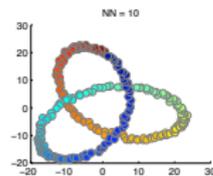
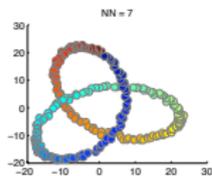
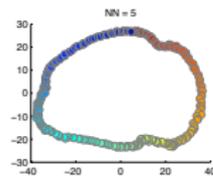
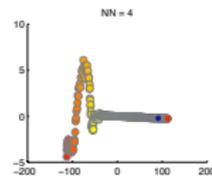
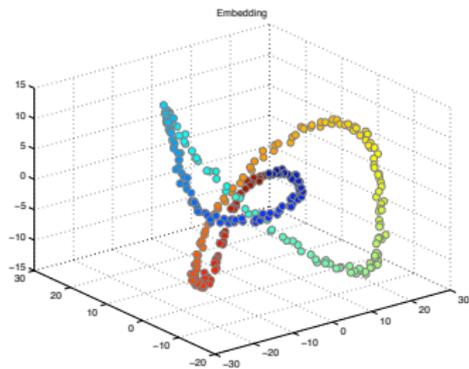
²/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example²



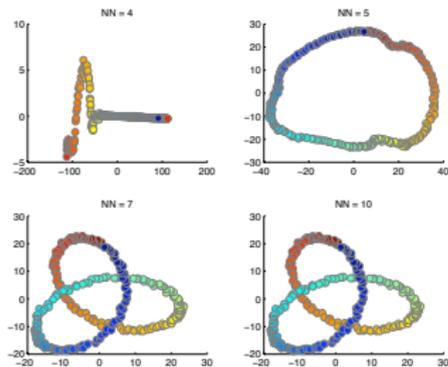
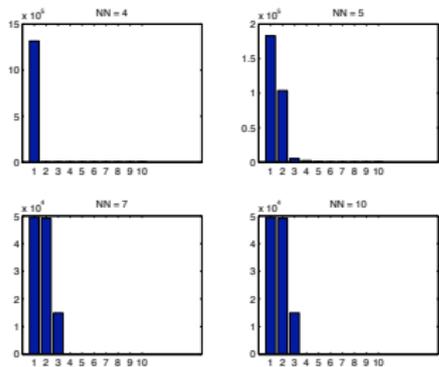
²/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example²



²/algos/mvu_embed.m

Maximum Variance Unfolding: Example²



²/algorithms/mvu_embed.m

Maximum Variance Unfolding: Summary

- MDS on optimised constrained Gram Matrix
- + Dimensionality through eigen spectra
- + Convex optimisation problem
- + Handles holes and non-convex manifolds
- Expensive

Locally Linear Embeddings

- *Roweis, Saul* - Science December 2000 (same issue as Isomap)
- Parametrise local geometry of data
- Extend local geometry globally

Locally Linear Embeddings

- Parametrise each point as a linear combination of its neighbors
- If each patch can be transformed by a translation, rotation and scaling to manifold
- \Rightarrow linear combination valid on manifold

Locally Linear Embeddings: Algorithm

- 1 Compute Proximity Graph
- 2 Compute Reconstruction Weights
- 3 Find low-dimensional embedding respecting weights

- Find weights in linear combination

$$\text{Minimize: } \epsilon = \sum_{i=1}^N \left\| \sum_{\mathbf{y}_j \in \{(\mathbf{y}_i, \mathbf{y}_j) \in W\}} \mathbf{w}_{ij} \mathbf{y}_j - \mathbf{y}_i \right\|_2^2$$

$$\text{Subject to: } \sum_{j \in \{(\mathbf{y}_i, \mathbf{y}_j) \in W\}} w_{ij} = 1$$

- Solution

$$\mathbf{w}_i = (\mathbf{N}^T \mathbf{N})^{-1} \left(\mathbf{N}^T \mathbf{y} - \frac{\mathbf{e}^T (\mathbf{N}^T \mathbf{N})^{-1} \mathbf{N}^T \mathbf{y} - 1}{\mathbf{e}^T (\mathbf{N}^T \mathbf{N})^{-1} \mathbf{e}} \right)$$

$$\mathbf{N} = [\mathbf{y}_{N(\mathbf{y}_i, 1)}, \dots, \mathbf{y}_{N(\mathbf{y}_i, K)}]^T$$

- Find low dimensional embedding \mathbf{X} respecting weights

$$\operatorname{argmin}_{\mathbf{X}} = \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{\mathbf{y}_j \in \{(\mathbf{y}_i, \mathbf{y}_j) \in W\}} w_{ij} \mathbf{x}_j \right\|_2^2$$

- Find \mathbf{X} that minimizes:

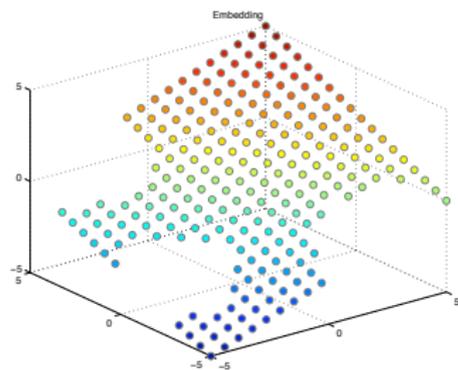
$$\mathbf{X}^T \underbrace{(\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})}_{\mathbf{M}} \mathbf{X}$$

- Objective function invariant to **scaling** and **translation**

$$\sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$$
$$\frac{1}{N-1} \mathbf{X}^T \mathbf{X} = \mathbf{I}$$

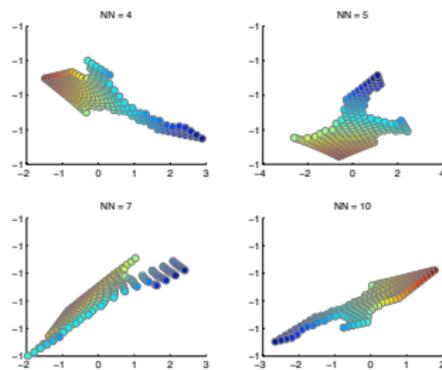
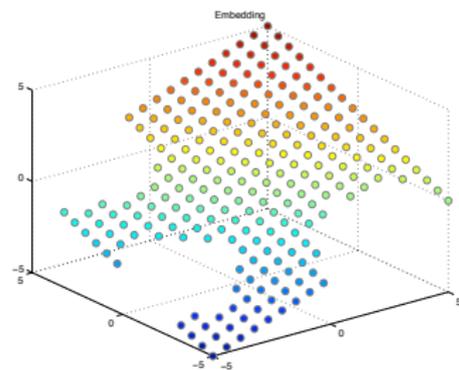
- Choose \mathbf{X} to be the *smallest* $\mathbf{d}+1$ eigenvectors of \mathbf{M}

Locally Linear Embeddings: Example³



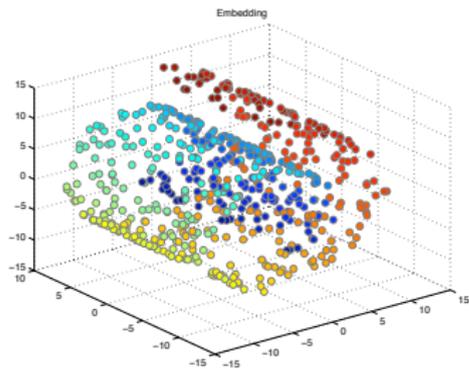
`3/algos/lle_embed.m`

Locally Linear Embeddings: Example³



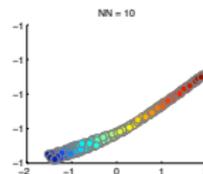
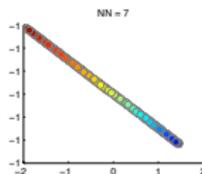
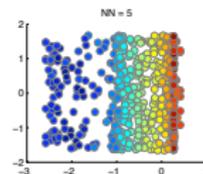
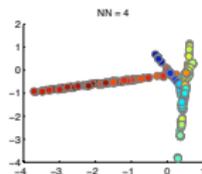
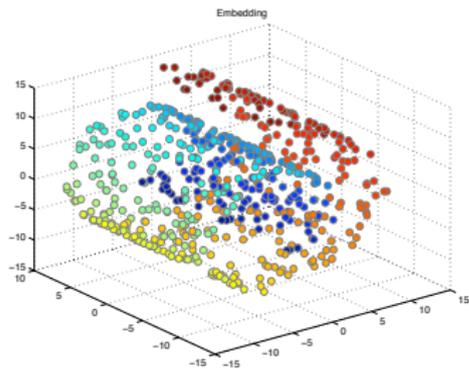
³/algos/lle_embed.m

Locally Linear Embeddings: Example³



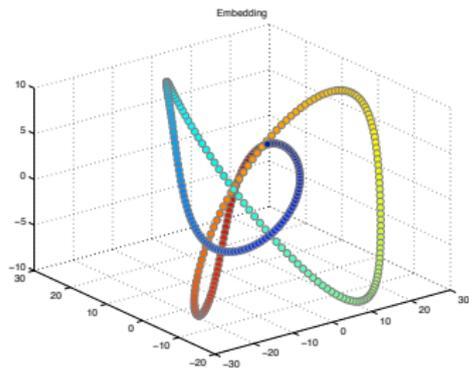
³/algorithms/lle_embed.m

Locally Linear Embeddings: Example³



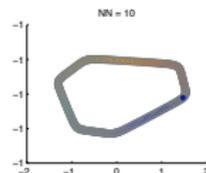
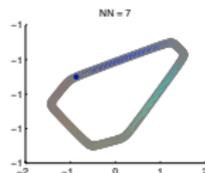
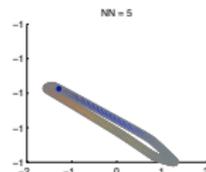
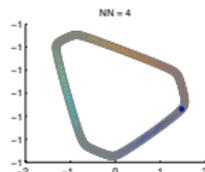
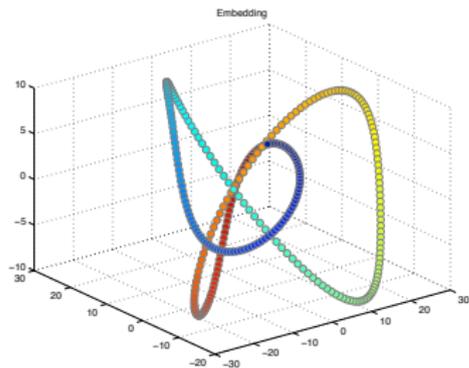
³/algos/lle_embed.m

Locally Linear Embeddings: Example³



³/algos/lle_embed.m

Locally Linear Embeddings: Example³



³/algorithms/lle_embed.m

Locally Linear Embeddings: Summary

- Unravel manifold by local parametrisation of each point
- + Solves a sparse eigenvalue problem
- + Finds bottom eigenvalues \Rightarrow Faster
- + handles holes and non-convex manifolds
 - Sensitive to non-uniform sampling
 - No indication of dimensionality
 - In practice hard to solve, (Matlabs eigensolver often fails)

- *Belkin, Niyogi* - NIPS 2001
- Find low dimensional embedding preserving locality
- Edgeweights correspond to locality measure

- Preserve “weighted” Locality

$$\begin{aligned} \operatorname{argmin}_{\mathbf{X}} &= \sum_{i=1}^N \sum_{j=1}^N \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \mathbf{W}_{ij} \\ & \quad (\mathbf{y}_i, \mathbf{y}_j) \in W \quad \begin{cases} \mathbf{w}_{ij} = e^{-\frac{\|\mathbf{y}_i - \mathbf{y}_j\|_2^2}{t}} \\ \mathbf{w}_{ij} = 1 \end{cases} \\ & \quad (\mathbf{y}_i, \mathbf{y}_j) \notin W \quad \mathbf{w}_{ij} = 0 \\ \operatorname{argmin}_{\mathbf{X}} &= \sum_{i=1}^N \sum_{j=1}^N \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \mathbf{W}_{ij} = \\ &= \{\mathbf{L} = \mathbf{D} - \mathbf{W}\} = \operatorname{trace}(\mathbf{X}^T \mathbf{L} \mathbf{X}) \end{aligned}$$

Laplacian Eigenmaps

- Trivial zero dimensional solution
- Remove scale invariance

$$\mathbf{x}^T \mathbf{D} \mathbf{1} = 0$$

$$\mathbf{x}^T \mathbf{D} \mathbf{x} = \mathbf{1}$$

- Objective

$$\begin{array}{ll} \operatorname{argmin}_{\mathbf{x}} & \operatorname{trace} \mathbf{X}^T \mathbf{L} \mathbf{X} \\ \text{subject to:} & \mathbf{x}^T \mathbf{D} \mathbf{1} = 0 \\ & \mathbf{x}^T \mathbf{D} \mathbf{x} = \mathbf{1} \end{array}$$

- Unconstrained solution given by the eigenvectors to \mathbf{L}
- Eigenvector corresponding to smallest eigenvalue $\lambda_N = 0$ corresponds to zero dimensional solution
- Constrained solution given by generalised eigenvalue problem

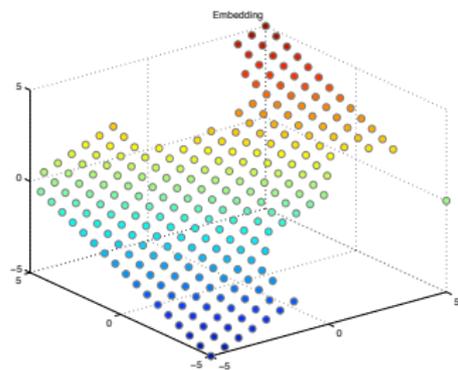
$$\mathbf{LX} = \mathbf{\Lambda DX}$$

- 1 Compute Proximity Graph
- 2 Complete Graph
- 3 Compute embedding from generalised eigenvalue problem

$$\mathbf{LX} = \mathbf{ADX}$$

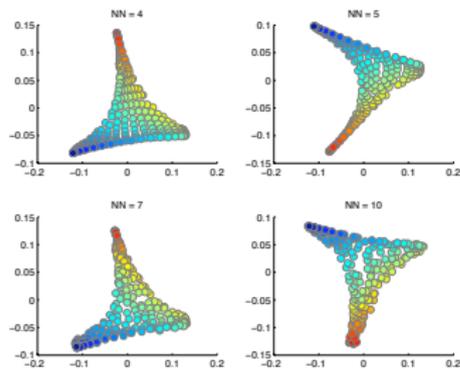
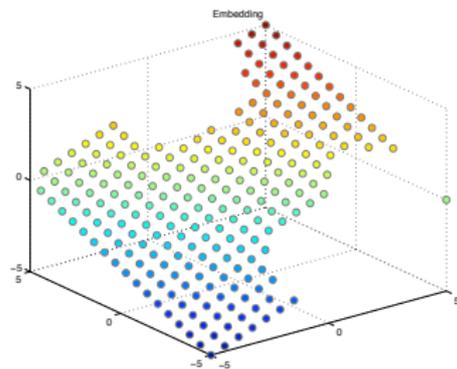
- 4 Embedding given by bottom $(d+1)$ generalised eigenvectors

Laplacian Eigenmaps: Example⁴



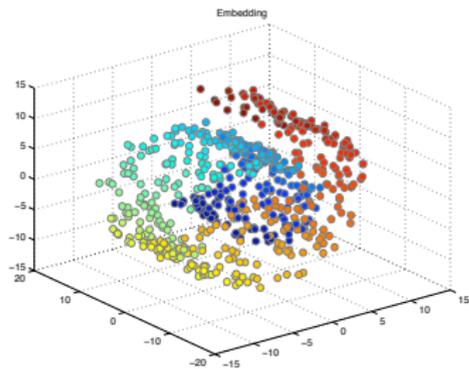
⁴ /algos/laplacian_embed.m

Laplacian Eigenmaps: Example⁴



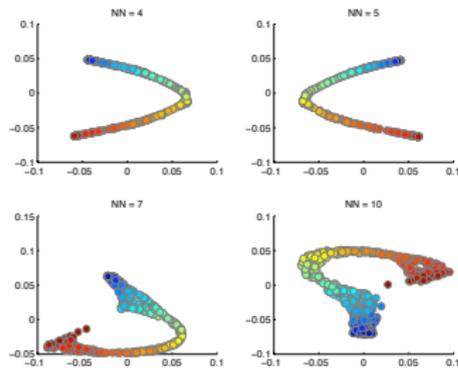
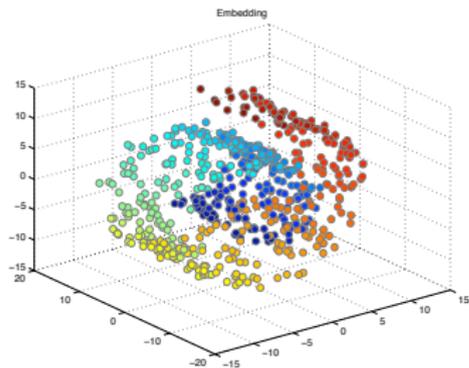
⁴/algorithms/laplacian_embed.m

Laplacian Eigenmaps: Example⁴



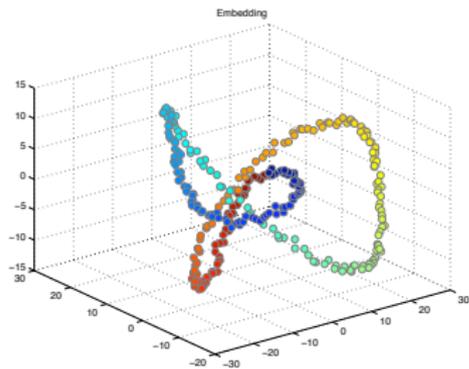
⁴/algos/laplacian_embed.m

Laplacian Eigenmaps: Example⁴



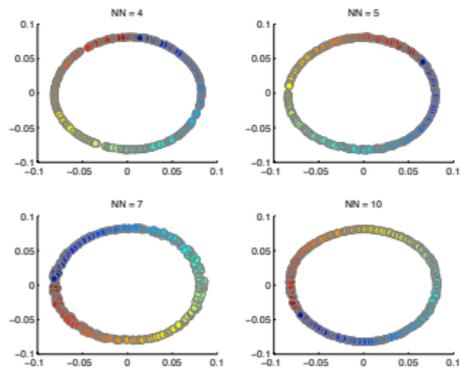
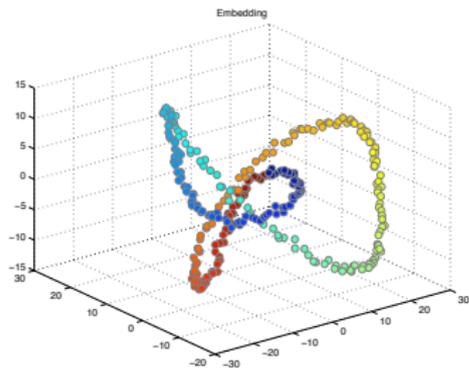
⁴/algos/laplacian_embed.m

Laplacian Eigenmaps: Example⁴



⁴ /algos/laplacian_embed.m

Laplacian Eigenmaps: Example⁴

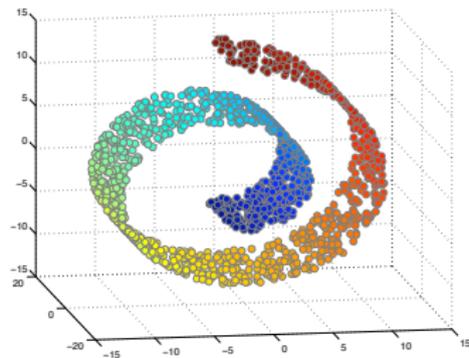


⁴/algorithms/laplacian_embed.m

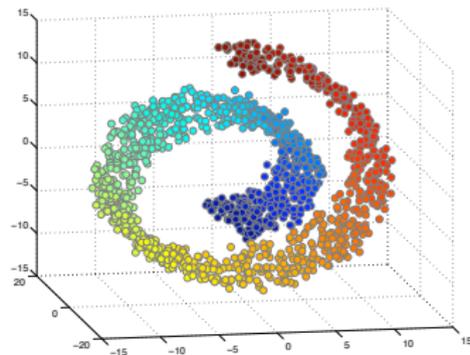
Laplacian Eigenmaps: Summary

- Unravels manifold by preserving locality
- + Finds bottom eigenvalues \Rightarrow Faster
- No indication of dimensionality

- Isomap and MVU non-linear extensions to MDS
- LLE preserves local parametrisation
- Laplacian Eigenmaps preserves locality

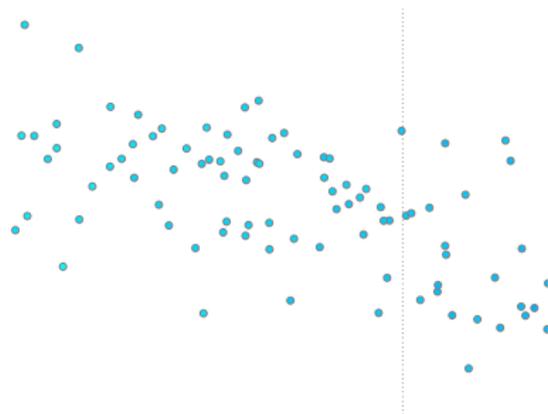
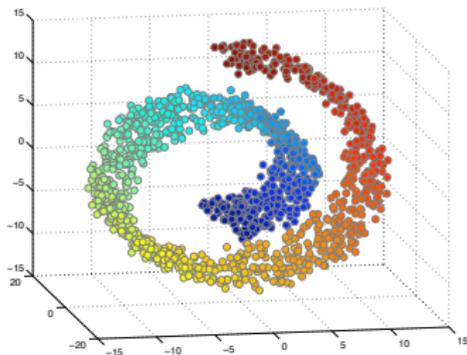


- Algorithms based on local assumption



- Algorithms based on local assumption
- Global noise viewed locally

Locality



- Algorithms based on local assumption
- Global noise viewed locally

- We have motivated the need for non-linear dimensionality reduction.
- Spectral approaches can achieve this, but they don't lead to probabilistic models.
- We are looking for a probabilistic approach to encoding the mapping.
- Next we will see how point based representations of the latent space can be used to achieve this.

Non Linear Probabilistic Methods I

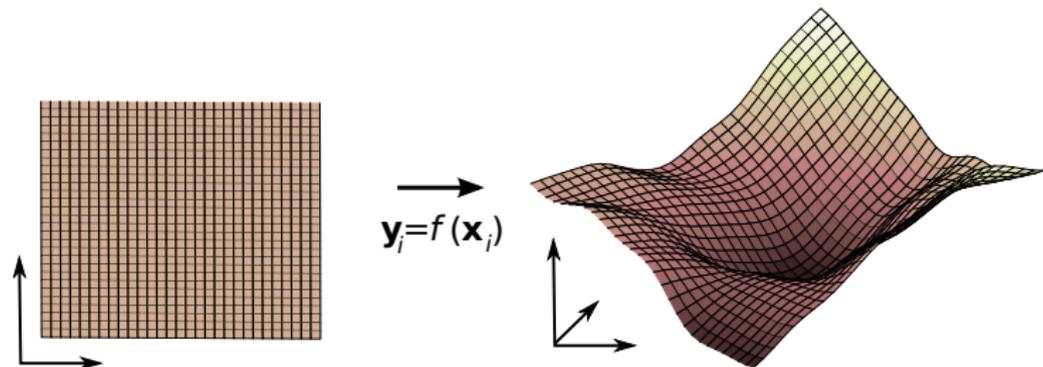


Figure: Mapping a two dimensional plane to a higher dimensional space in a non-linear way.

Difficulty for Probabilistic Approaches

- Propagate a probability distribution through a non-linear mapping.
- Normalisation of distribution becomes intractable.

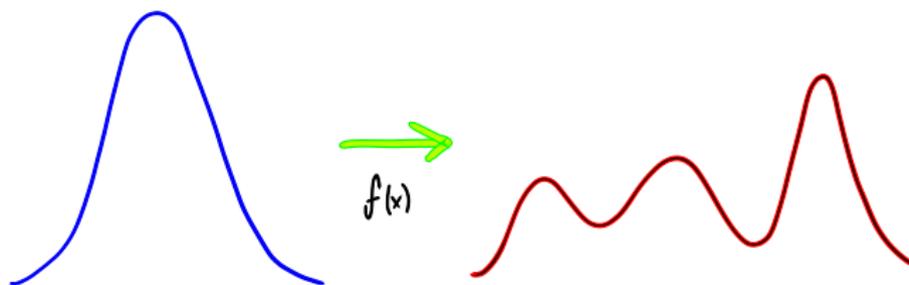


Figure: Gaussian distribution propagated through a non-linear mapping.

Sampling Approach

- Proposed as Density Networks (MacKay, 1995)
- Likelihood is a Gaussian with non-linear mapping from latent space to data space for the mean

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{i=1}^N \prod_{j=1}^D \mathcal{N}(y_{i,j} | f_j(\mathbf{x}_{i,:}; \boldsymbol{\theta}), \sigma^2)$$

$$p(\mathbf{X}) = \mathcal{N}(\mathbf{x}_{i,:} | \mathbf{0}, \mathbf{I})$$

- Take the mapping to be e.g. a multi-layer perceptron.
- Key idea: share same samples for all data points $\hat{\mathbf{X}}_n = \hat{\mathbf{X}} = \{\hat{\mathbf{x}}_{k,:}\}_{k=1}^M$.
- Saves computation — compute the mapping M times instead of MN

Mapping of Points

- Mapping points to higher dimensions is easy.

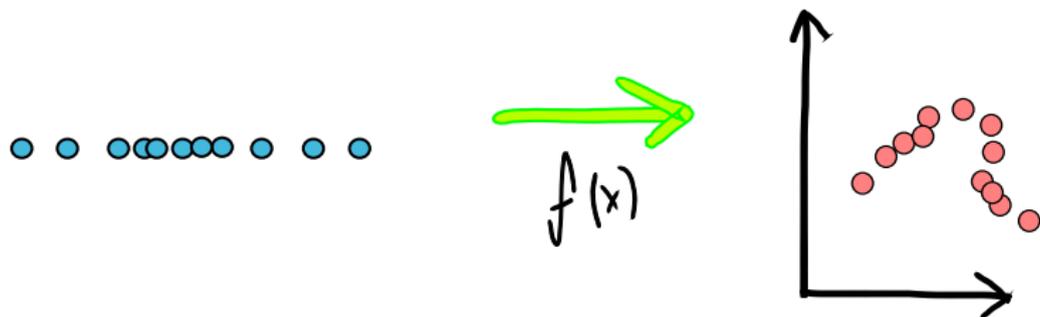


Figure: One dimensional Gaussian mapped to two dimensions.

Mapping of Points

- Mapping points to higher dimensions is easy.

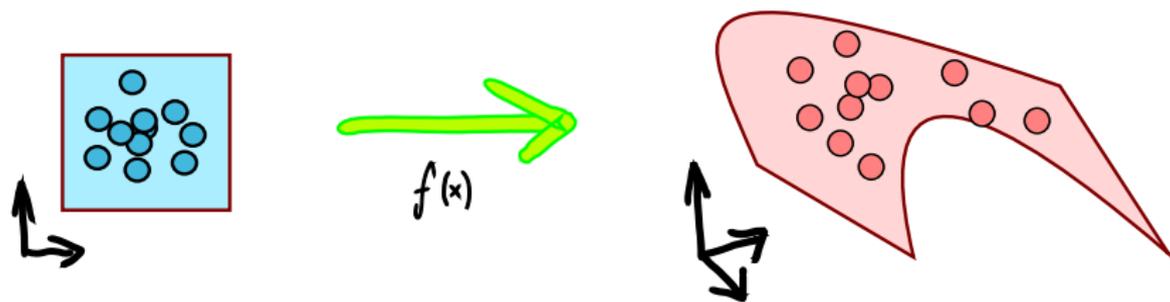


Figure: Two dimensional Gaussian mapped to three dimensions.

Sample approximation to log likelihood:

$$\log p(\mathbf{Y}|\boldsymbol{\theta}) = \sum_{i=1}^N \log \frac{1}{M} \sum_{k=1}^M p(\mathbf{y}_{i,:}|\boldsymbol{\theta}, \hat{\mathbf{x}}_{k,:})$$

so we have

$$\frac{d}{d\boldsymbol{\theta}} \log p(\mathbf{y}_{i,:}|\boldsymbol{\theta}) = \sum_{k=1}^M \frac{p(\mathbf{y}_{i,:}|\boldsymbol{\theta}, \hat{\mathbf{x}}_{k,:})}{\sum_{m=1}^M p(\mathbf{y}_{i,:}|\boldsymbol{\theta}, \hat{\mathbf{x}}_{m,:})} \frac{d}{d\boldsymbol{\theta}} \log p(\mathbf{y}_{i,:}|\boldsymbol{\theta}, \hat{\mathbf{x}}_{k,:})$$

$$\frac{d}{d\boldsymbol{\theta}} \log p(\mathbf{y}_{i,:}|\boldsymbol{\theta}) = \sum_{k=1}^M \hat{\pi}_{i,k} \frac{d}{d\boldsymbol{\theta}} \log p(\mathbf{y}_{i,:}|\boldsymbol{\theta}, \hat{\mathbf{x}}_{k,:})$$

Note: $\hat{\pi}_{i,k}$ look a bit like the posterior over component k for data point i .

- Use gradient based optimisation to find the mapping.

Generative Topographic Mapping

- Generative Topographic Mapping (GTM) (Bishop et al., 1998a)
- Key idea: Lay points out on a *grid*.
 - Constrained mixture of Gaussians.

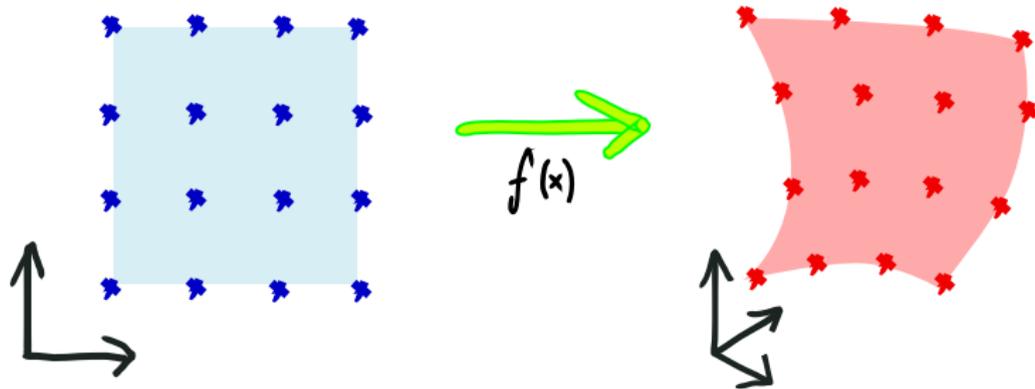


Figure: One dimensional Gaussian mapped to two dimensions.

- Prior distribution is a mixture model in a latent space.

$$p(\mathbf{X}) = \prod_{i=1}^N p(\mathbf{x}_{i,:})$$

$$p(\mathbf{x}_{i,:}) = \frac{1}{M} \sum_{k=1}^M \delta(\mathbf{x}_{i,:} - \hat{\mathbf{x}}_{k,:})$$

- The $\hat{\mathbf{x}}_{k,:}$ are laid out on a regular grid.

Mapping and E-Step

- Likelihood is a Gaussian with non-linear mapping from latent space to data space for the mean

$$p(\mathbf{Y}|\mathbf{X}, \theta) = \prod_{i=1}^N \prod_{j=1}^D \mathcal{N}(y_{i,j} | f_j(\mathbf{x}_{i,:}; \mathbf{W}, l), \sigma^2)$$

In the original paper (Bishop et al., 1998b) an RBF network was suggested,

- In the E-step, posterior distribution over k is given by

$$\hat{\pi}_{i,k} = \frac{\prod_{j=1}^D \mathcal{N}(y_{i,j} | f_j(\hat{\mathbf{x}}_k; \mathbf{W}, l), \sigma^2)}{\sum_{m=1}^M \prod_{j=1}^D \mathcal{N}(y_{i,j} | f_j(\hat{\mathbf{x}}_m; \mathbf{W}, l), \sigma^2)}$$

sometimes called the “responsibility of component k for data point i ”.

Likelihood Optimisation

- We then maximise the lower bound on the log likelihood,

$$\log p(\mathbf{y}_{i,:}|\boldsymbol{\theta}) \geq \langle \log p(\mathbf{y}_{i,:}, \hat{\mathbf{x}}_{k,:}|\boldsymbol{\theta}) \rangle_{q(k)} - \langle \log q(k) \rangle_{q(k)},$$

- Free energy part of bound

$$\langle \log p(\mathbf{y}_{i,:}, \hat{\mathbf{x}}_{k,:}|\boldsymbol{\theta}) \rangle = \sum_{k=1}^M \hat{\pi}_{i,k} \log p(\mathbf{y}_{i,:}|\hat{\mathbf{x}}_{k,:}, \boldsymbol{\theta}) + \text{const}$$

- When optimising parameters in EM, we ignore dependence of $\hat{\pi}_{i,k}$ on parameters. So we have

$$\frac{d}{d\boldsymbol{\theta}} \langle \log p(\mathbf{y}_{i,:}, \hat{\mathbf{x}}_{k,:}|\boldsymbol{\theta}) \rangle = \sum_{k=1}^M \hat{\pi}_{i,k} \frac{d}{d\boldsymbol{\theta}} \log p(\mathbf{y}_{i,:}|\hat{\mathbf{x}}_{k,:}, \boldsymbol{\theta})$$

which is very similar to density network result!

- Interpretation of posterior is slightly different.

Stick Man Data

- $N = 55$ frames of motion capture.
- xyz locations of 34 points on the body.
- $D = 102$ dimensional data.
- “Run 1” available from http://accad.osu.edu/research/mocap/mocap_data.htm.

Changing



Angle



of Run



demStickDnet1

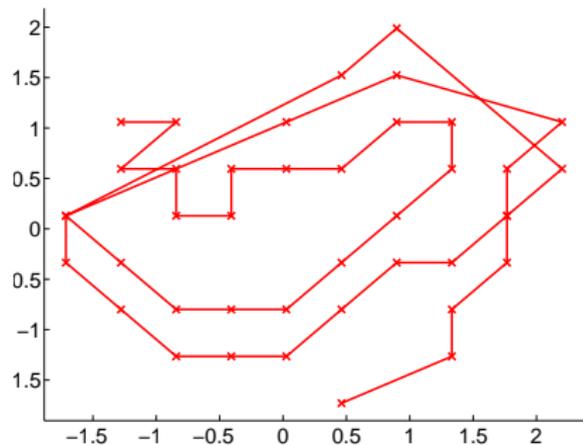


Figure: Stick man data visualised with the GTM using an RBF network with 10×10 points in the grid.

demStickDnet2

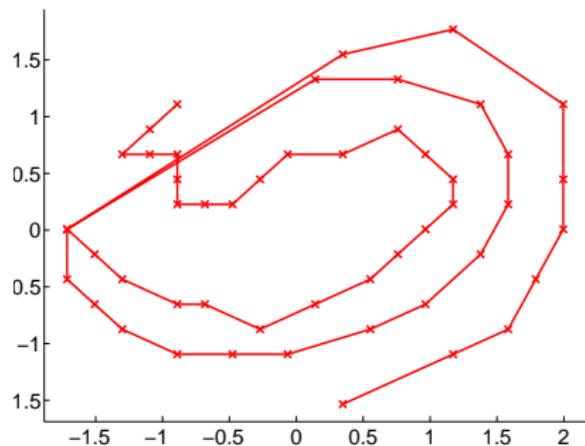


Figure: Stick man data visualised with the GTM using an RBF network with 20×20 points in the grid.

Bubblewrap Effect

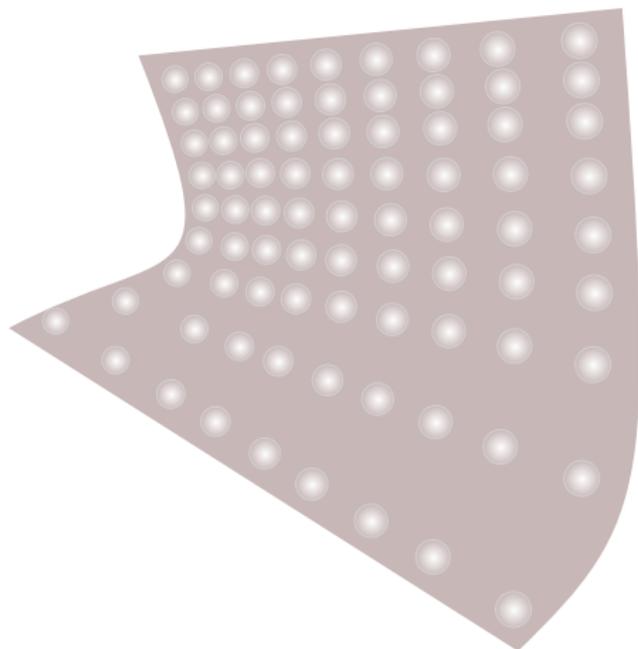
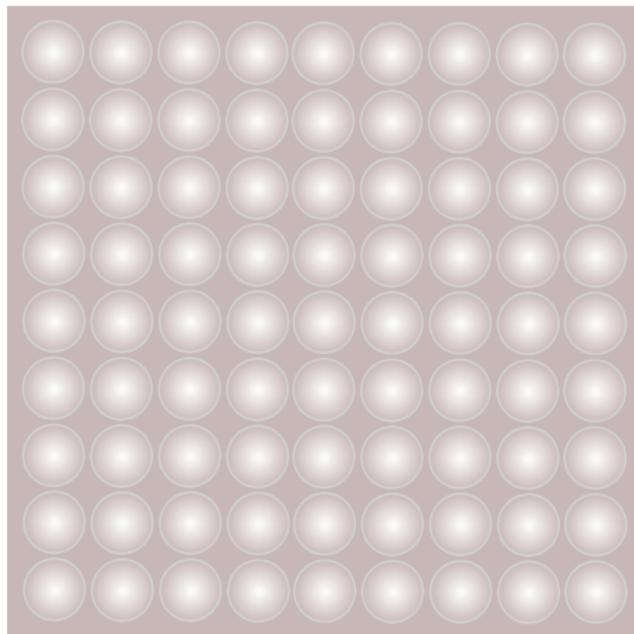


Figure: The manifold is more like bubblewrap than a piece of paper.

Effect of Separated Means

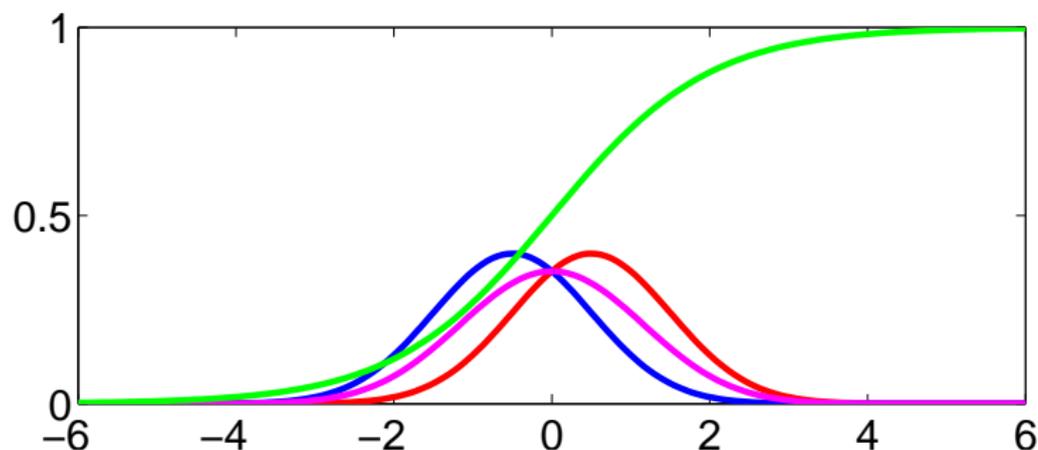


Figure: As Gaussians become further apart the posterior probability becomes more abrupt. 1 standard deviations apart.

Effect of Separated Means

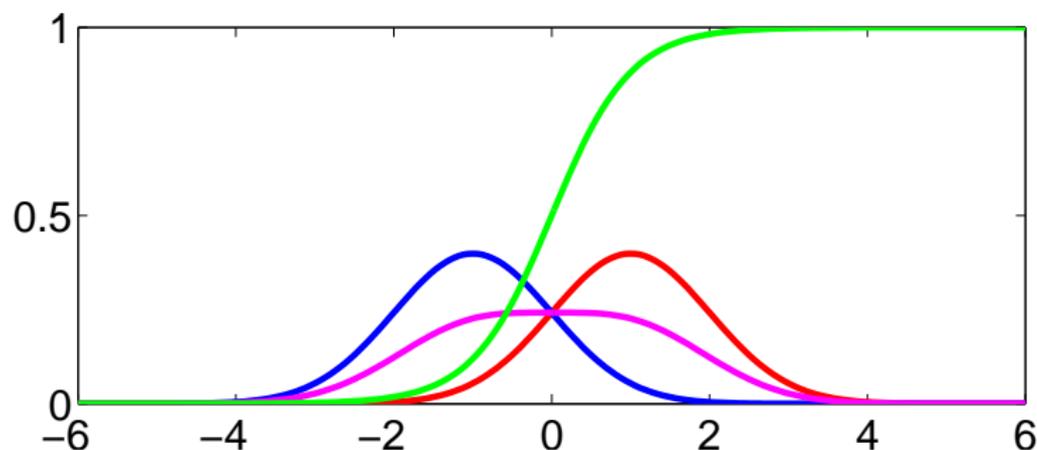


Figure: As Gaussians become further apart the posterior probability becomes more abrupt. 2 standard deviations apart.

Effect of Separated Means

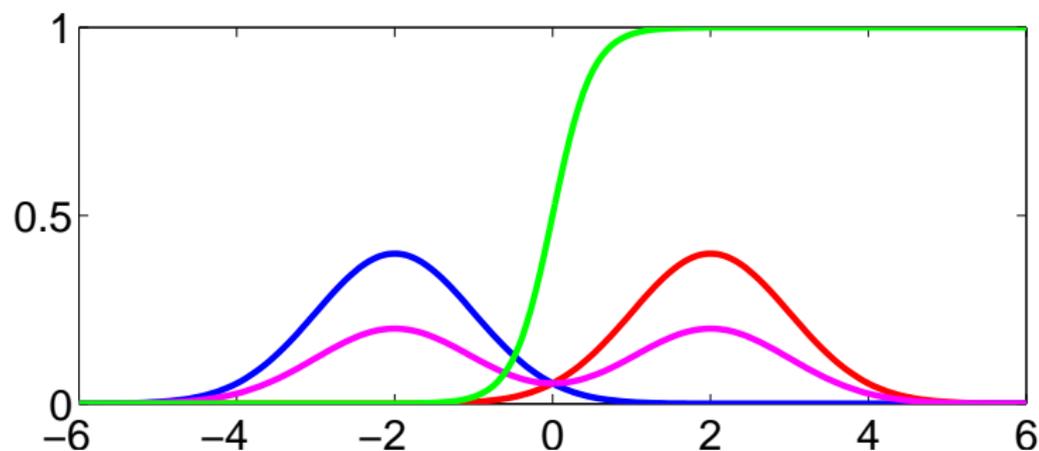


Figure: As Gaussians become further apart the posterior probability becomes more abrupt. 4 standard deviations apart.

Effect of Separated Means

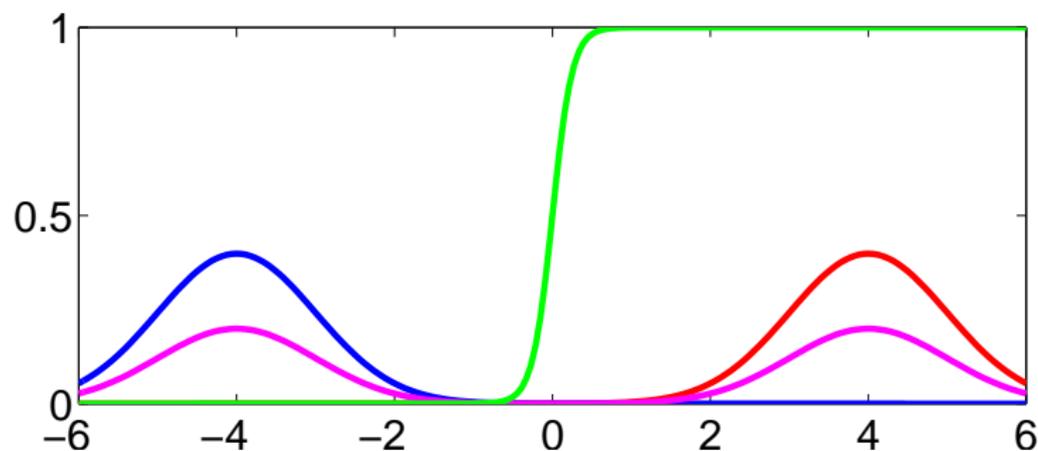


Figure: As Gaussians become further apart the posterior probability becomes more abrupt. 8 standard deviations apart.

Effect of Separated Means

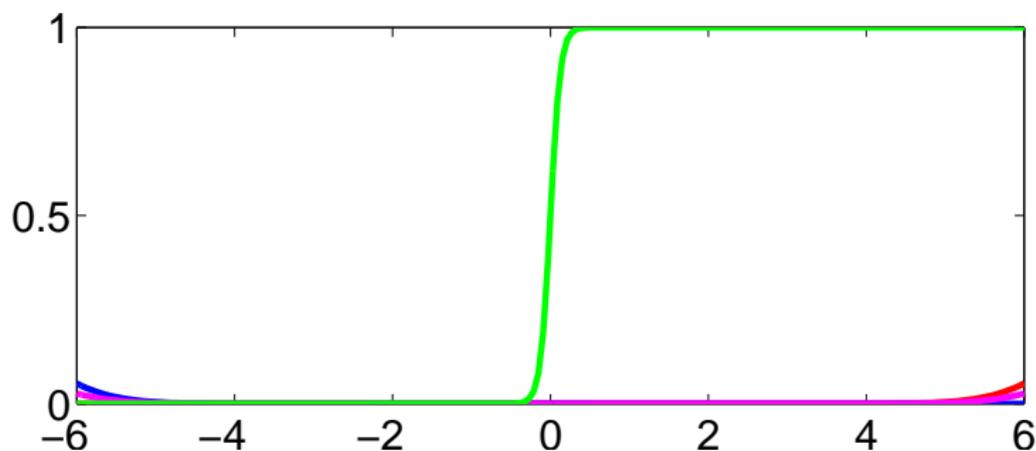


Figure: As Gaussians become further apart the posterior probability becomes more abrupt. 16 standard deviations apart.

Equivalence of GTM and Density Networks

- GTM and Density Networks have the same origin. (Bishop et al. 1996; McKay, 1995).
- In original Density Networks paper MacKay suggested Importance Sampling (MacKay, 1995).
- Early work on GTM also used importance sampling.
- Main innovation in GTM was to lay points out on a grid (inspired by Self Organizing Maps (Kohonen, 2001)).

- We have explored two point based approaches to dimensionality reduction.
- Approaches seem to generalise well even when dimensions of data is greater than number of points.
- Both approaches are difficult to extend to higher dimensional latent spaces
 - number of samples/centres required increases exponentially with dimension.
- Next we will explore a different probabilistic interpretation of PCA and extend that to non-linear models.

Probabilistic PCA

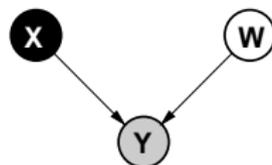
- We have seen that PCA has a probabilistic interpretation (Tipping and Bishop, 1999b) .
- It is difficult to 'non-linearise' directly.
- GTM and Density Networks are an attempt to do so.

Dual Probabilistic PCA

- There is an alternative probabilistic interpretation of PCA (Lawrence, 2005) .
- This interpretation can be made non-linear.
- The result is non-linear probabilistic PCA.

Dual Probabilistic PCA

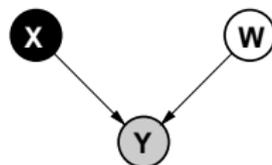
- Define *linear-Gaussian relationship* between latent variables and data.
 - **Novel** Latent variable approach:



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA

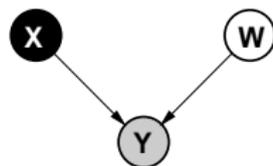
- Define *linear-Gaussian relationship* between latent variables and data.
 - **Novel** Latent variable approach:
 - Define Gaussian prior over *parameters*, \mathbf{W} .



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(y_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
 - **Novel** Latent variable approach:
 - Define Gaussian prior over *parameters*, \mathbf{W} .
 - Integrate out *parameters*.

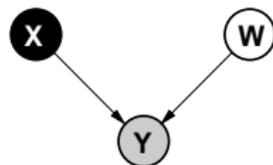


$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{i=1}^D \mathcal{N}(\mathbf{w}_{i,:} | \mathbf{0}, \mathbf{I})$$

Dual Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
 - **Novel** Latent variable approach:
 - Define Gaussian prior over *parameters*, \mathbf{W} .
 - Integrate out *parameters*.

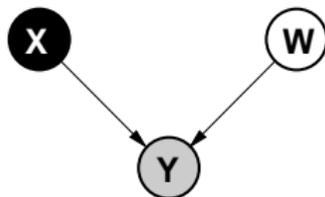


$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^N \mathcal{N}(y_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{i=1}^D \mathcal{N}(\mathbf{w}_{i,:} | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D \mathcal{N}(y_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA Max. Likelihood Soln (Lawrence, 2004)



$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D \mathcal{N}(y_{:j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA Max. Likelihood Soln (Lawrence, 2004)

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D \mathcal{N}(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{K}), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

$$\log p(\mathbf{Y}|\mathbf{X}) = -\frac{D}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^T) + \text{const.}$$

If \mathbf{U}'_q are first q principal eigenvectors of $D^{-1}\mathbf{Y}\mathbf{Y}^T$ and the corresponding eigenvalues are Λ_q ,

$$\mathbf{X} = \mathbf{U}'_q \mathbf{L} \mathbf{R}^T, \quad \mathbf{L} = (\Lambda_q - \sigma^2 \mathbf{I})^{\frac{1}{2}}$$

where \mathbf{R} is an arbitrary rotation matrix.

Probabilistic PCA Max. Likelihood Soln (Tipping and Bishop, 1999b)

$$p(\mathbf{Y}|\mathbf{W}) = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{C}), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$$

$$\log p(\mathbf{Y}|\mathbf{W}) = -\frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \text{tr}(\mathbf{C}^{-1}\mathbf{Y}^T\mathbf{Y}) + \text{const.}$$

If \mathbf{U}_q are first q principal eigenvectors of $N^{-1}\mathbf{Y}^T\mathbf{Y}$ and the corresponding eigenvalues are Λ_q ,

$$\mathbf{W} = \mathbf{U}_q\mathbf{L}\mathbf{R}^T, \quad \mathbf{L} = (\Lambda_q - \sigma^2\mathbf{I})^{\frac{1}{2}}$$

where \mathbf{R} is an arbitrary rotation matrix.

The Eigenvalue Problems are equivalent

- Solution for Probabilistic PCA (solves for the mapping)

$$\mathbf{Y}^T \mathbf{Y} \mathbf{U}_q = \mathbf{U}_q \Lambda_q \quad \mathbf{W} = \mathbf{U}_q \mathbf{L} \mathbf{V}^T$$

- Solution for Dual Probabilistic PCA (solves for the latent positions)

$$\mathbf{Y} \mathbf{Y}^T \mathbf{U}'_q = \mathbf{U}'_q \Lambda_q \quad \mathbf{X} = \mathbf{U}'_q \mathbf{L} \mathbf{V}^T$$

- Equivalence is from

$$\mathbf{U}_q = \mathbf{Y}^T \mathbf{U}'_q \Lambda_q^{-\frac{1}{2}}$$

Prior for Functions

- Probability Distribution over Functions
- Functions are infinite dimensional.
 - Prior distribution over *instantiations* of the function: finite dimensional objects.
 - Can prove by induction that GP is 'consistent'.
- Mean and Covariance Functions
- Instead of mean and covariance matrix, GP is defined by mean function and covariance function.
 - Mean function often taken to be zero or constant.
 - Covariance function must be *positive definite*.
 - Class of valid covariance functions is the same as the class of *Mercer kernels*.

Zero mean Gaussian Process

- A (zero mean) Gaussian process likelihood is of the form

$$p(\mathbf{y}|\mathbf{X}) = N(\mathbf{y}|\mathbf{0}, \mathbf{K}),$$

where \mathbf{K} is the covariance function or *kernel*.

- The *linear kernel* with noise has the form

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

- Priors over non-linear functions are also possible.
 - To see what functions look like, we can sample from the prior process.

demCovFuncSample

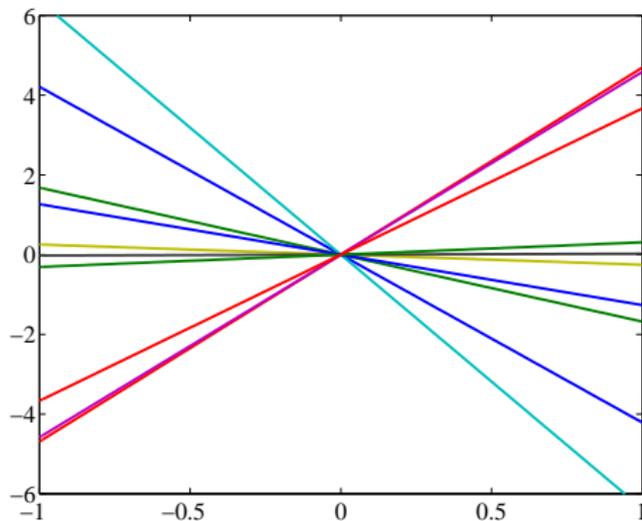


Figure: linear kernel, $\mathbf{K} = \mathbf{X}\mathbf{X}^T$

demCovFuncSample

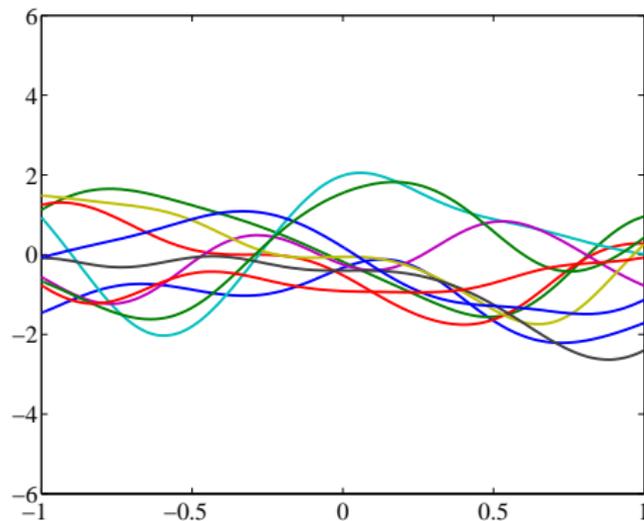


Figure: RBF kernel with $\gamma = 10$, $\alpha = 1$

demCovFuncSample

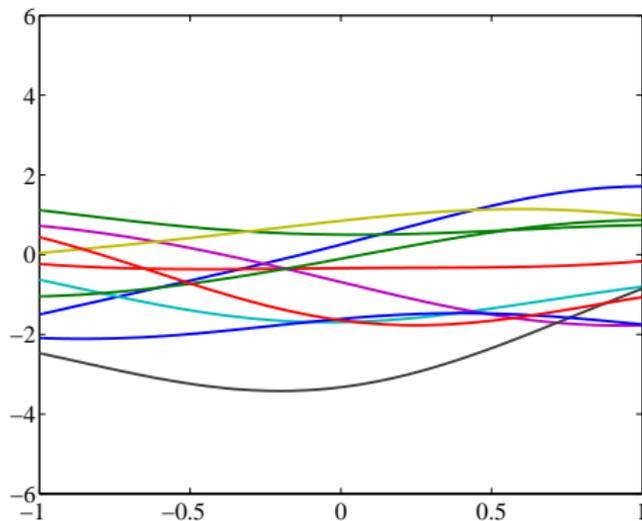


Figure: RBF kernel with $l = 1$, $\alpha = 1$

demCovFuncSample

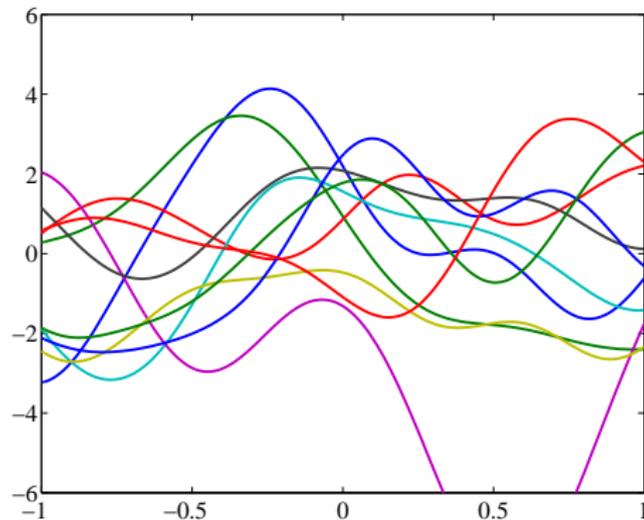


Figure: RBF kernel with $l = 0.3$, $\alpha = 4$

demCovFuncSample

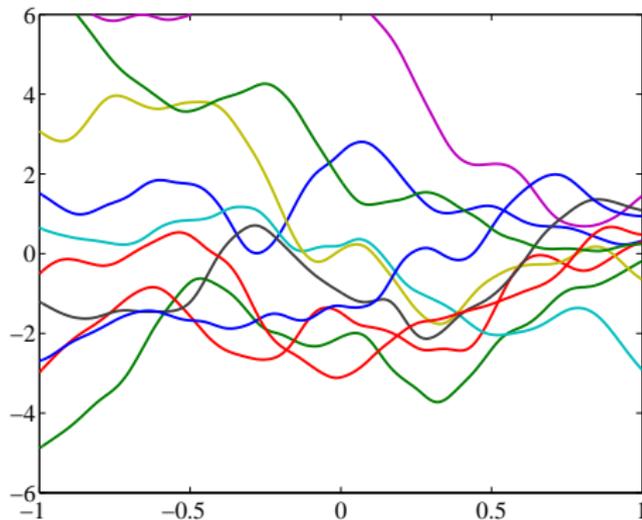


Figure: MLP kernel with $\alpha = 8$, $w = 100$ and $b = 100$

demCovFuncSample

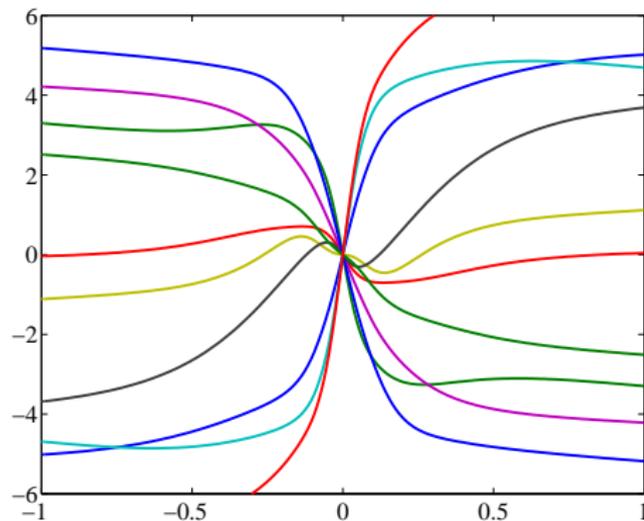


Figure: MLP kernel with $\alpha = 8$, $b = 0$ and $w = 100$

demCovFuncSample

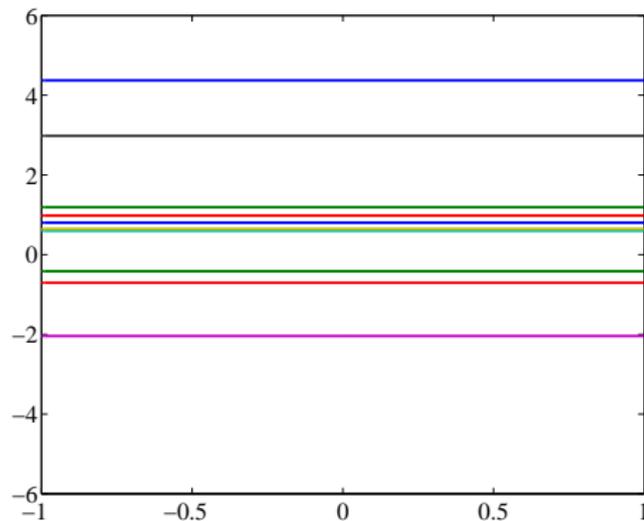


Figure: bias kernel with $\alpha = 1$ and

Covariance Samples

demCovFuncSample

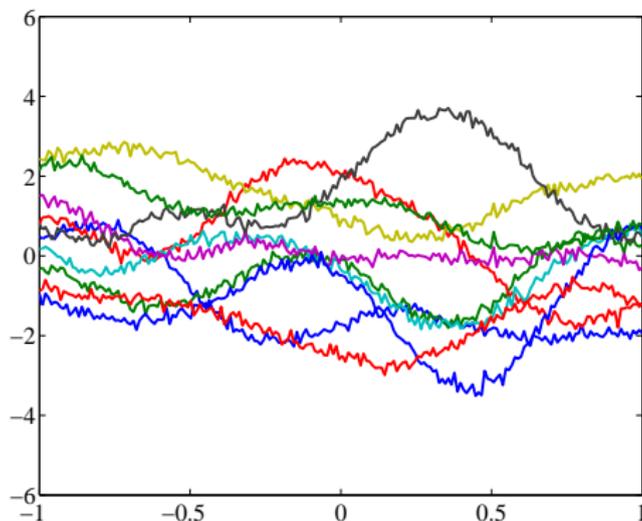


Figure: summed combination of: RBF kernel, $\alpha = 1$, $l = 0.3$; bias kernel, $\alpha = 1$; and white noise kernel, $\beta = 100$

Posterior Distribution over Functions

- Gaussian processes are often used for regression.
- We are given a known inputs \mathbf{X} and targets \mathbf{Y} .
- We assume a prior distribution over functions by selecting a kernel.
- Combine the prior with data to get a *posterior* distribution over functions.

Gaussian Process Regression

demRegression

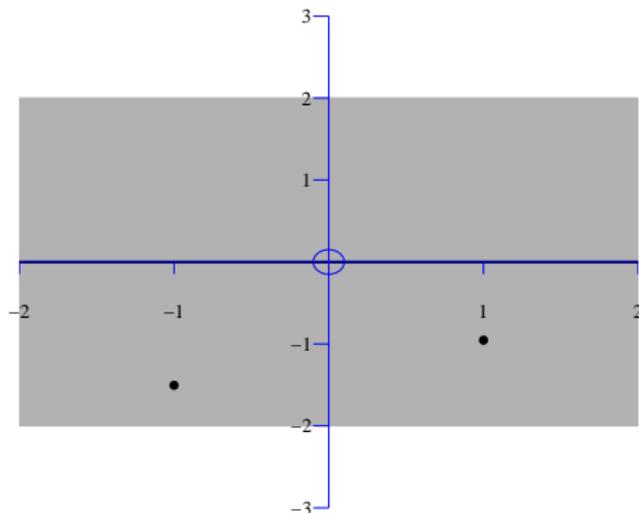


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

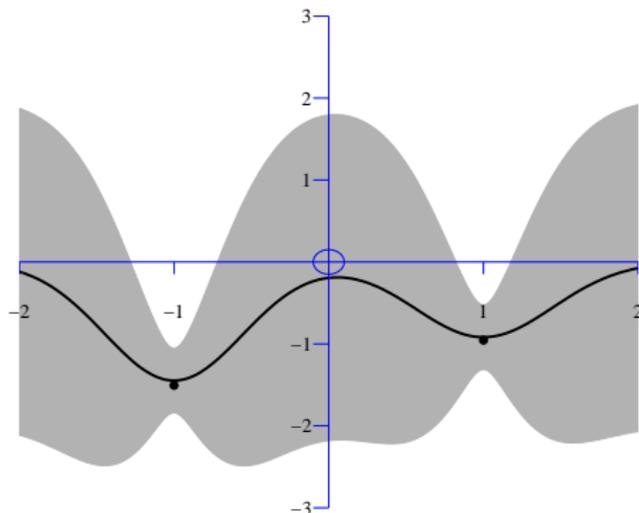


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

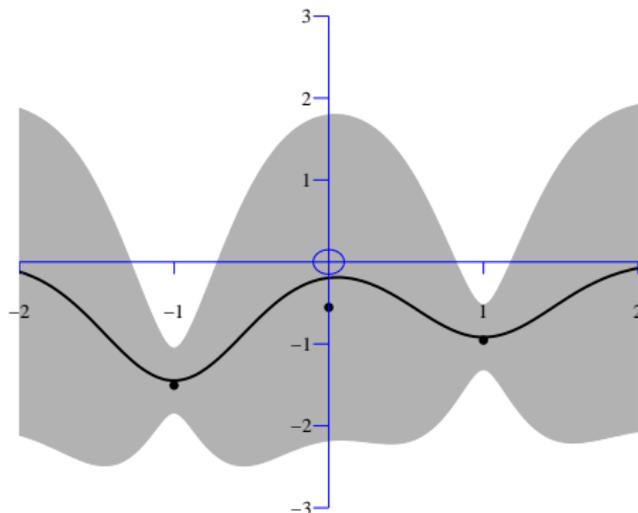


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

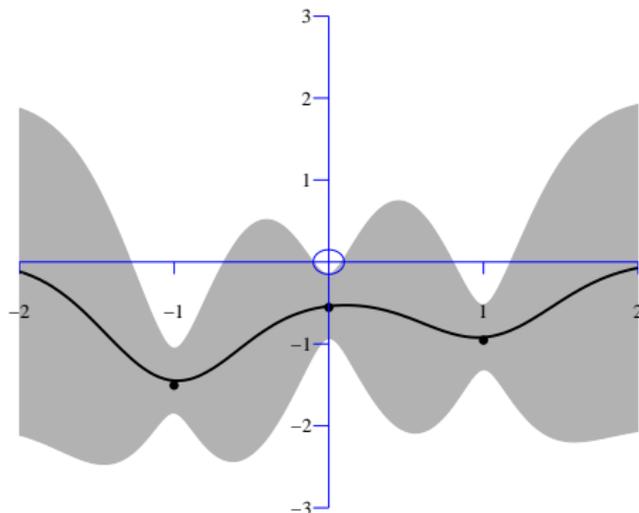


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

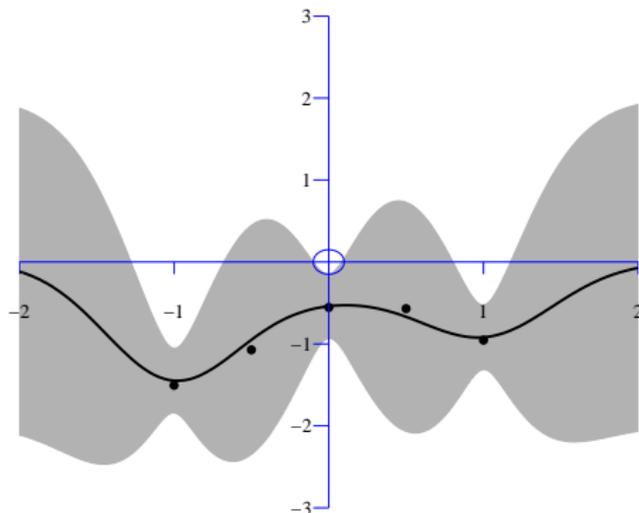


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

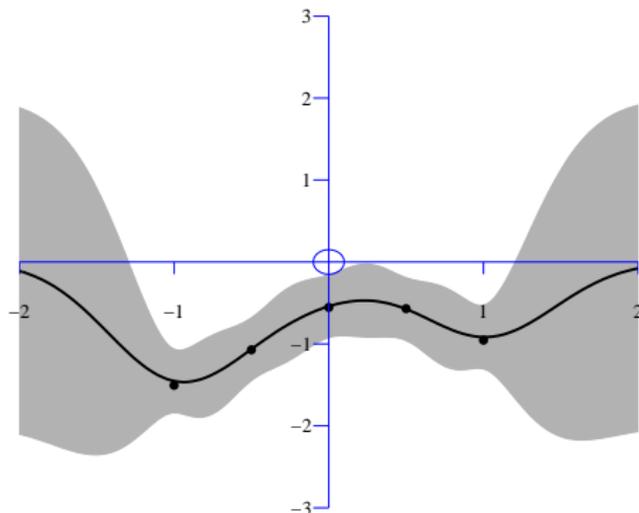


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

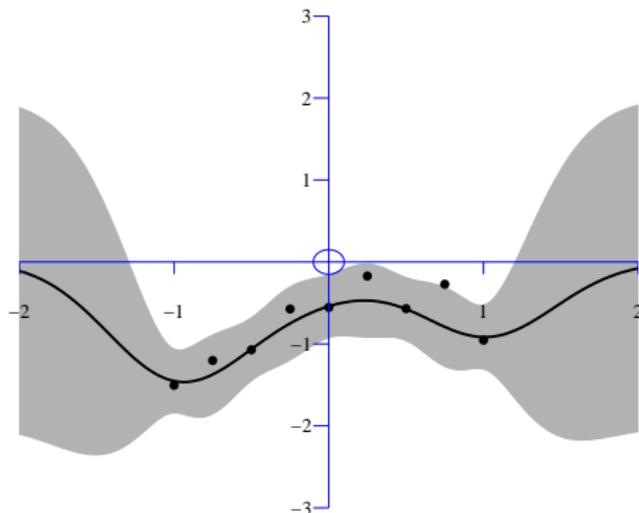


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

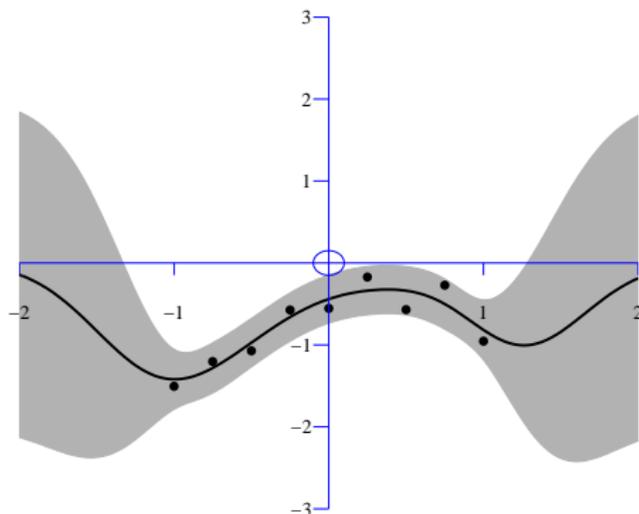
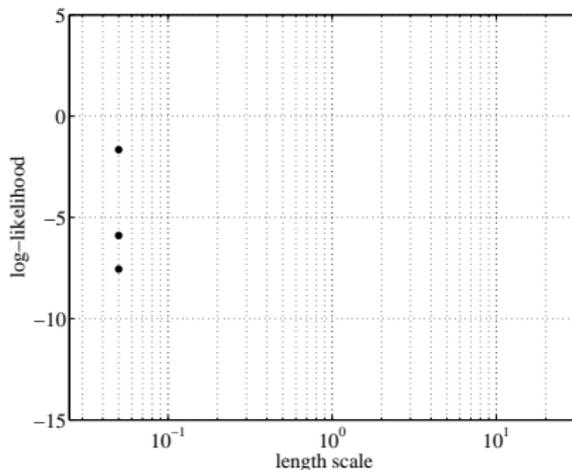
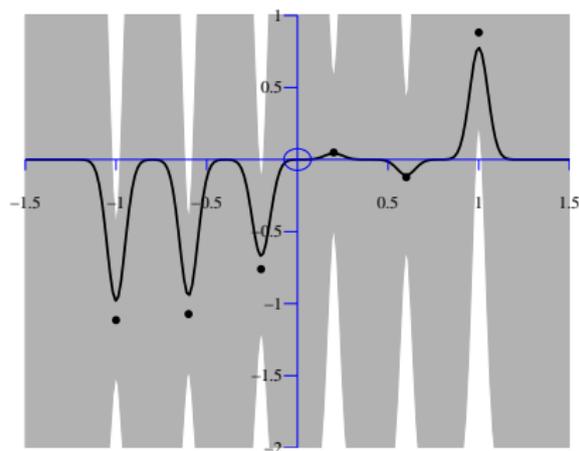


Figure: Examples include WiFi localization, C14 calibration curve.

Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

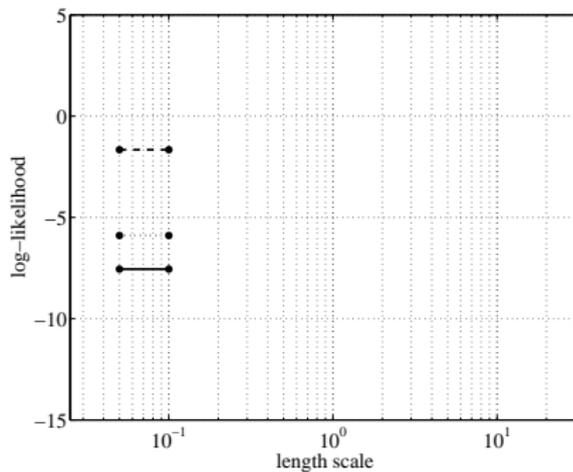
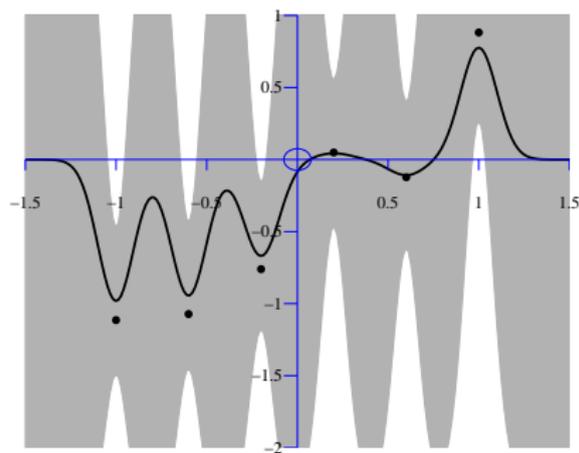
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

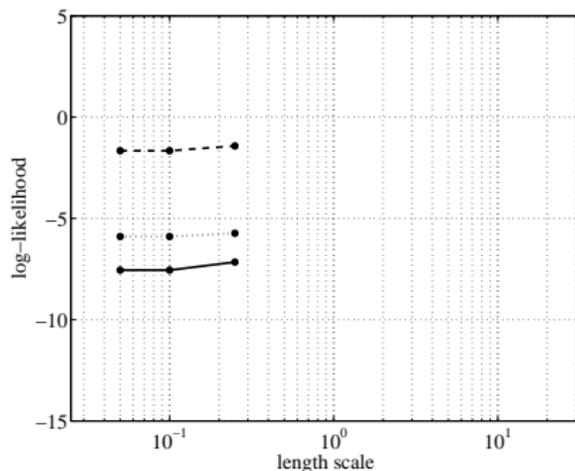
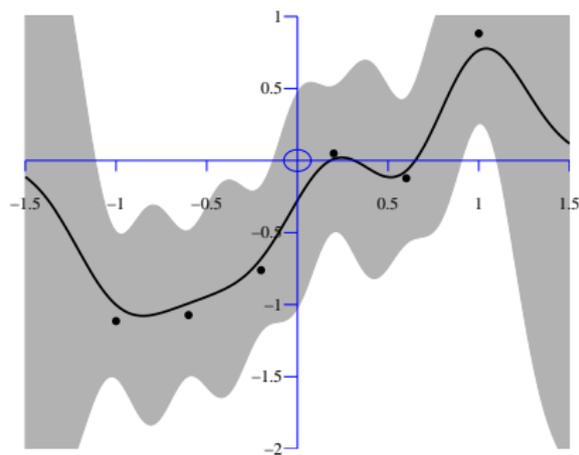
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

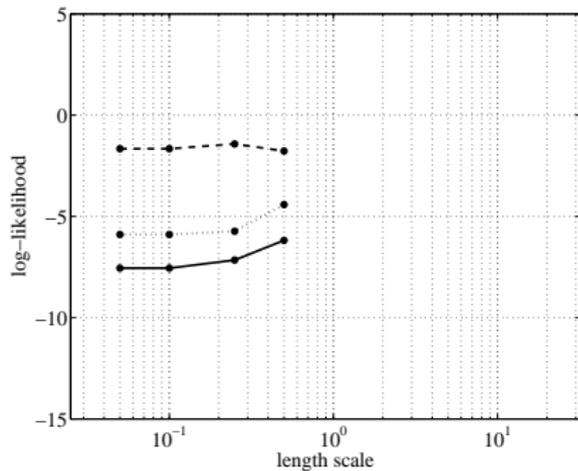
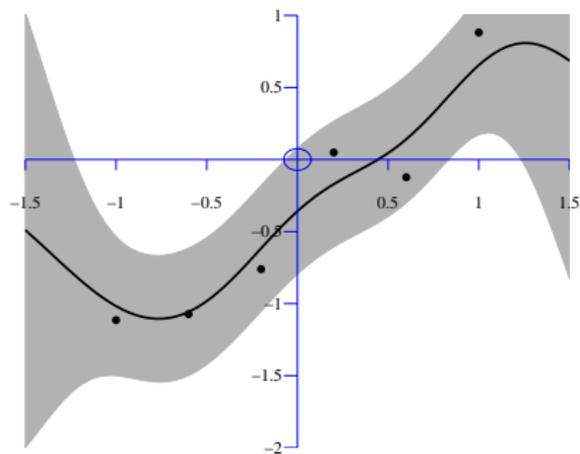
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

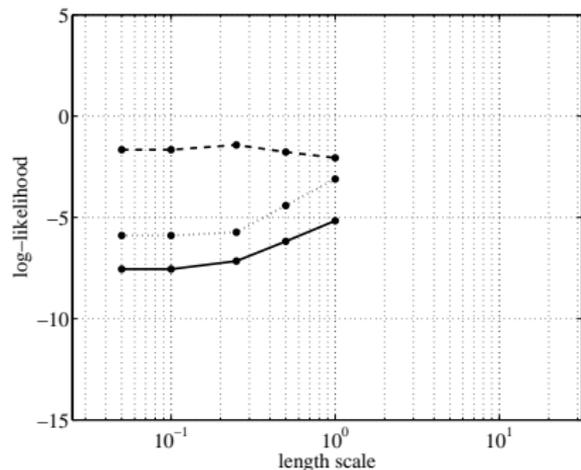
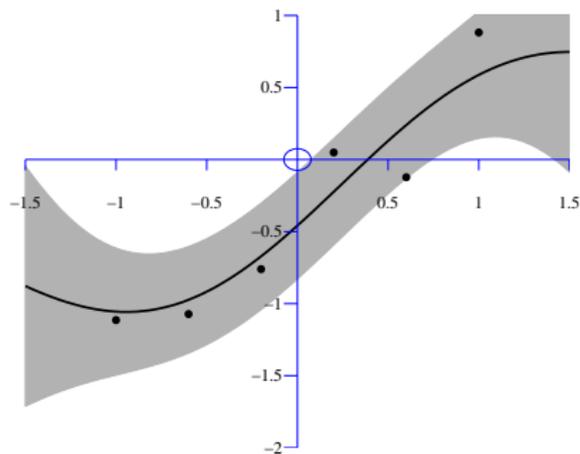
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

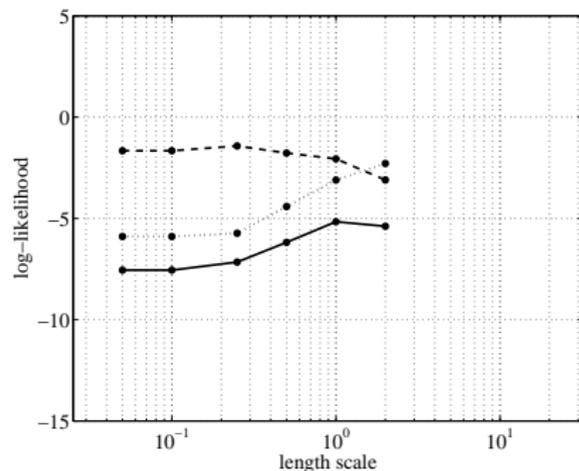
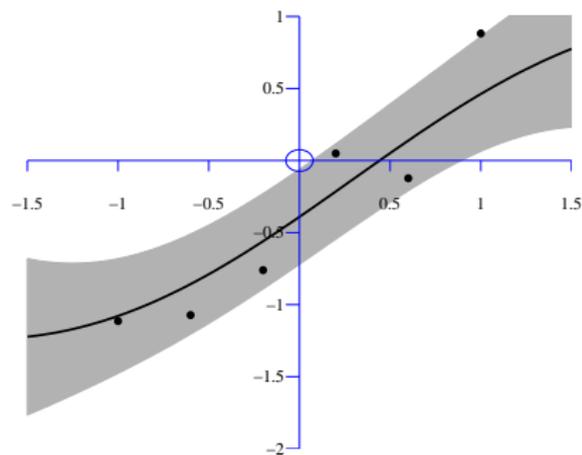
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

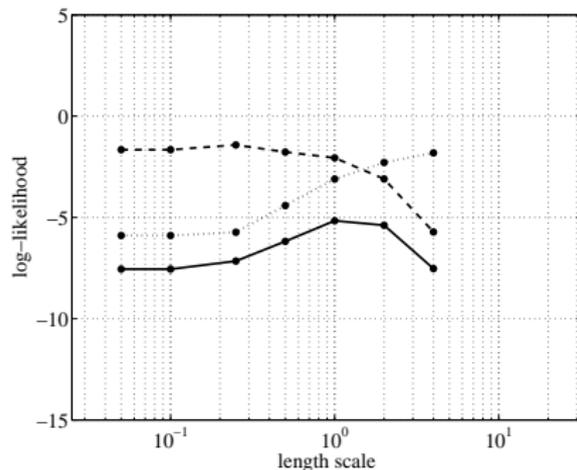
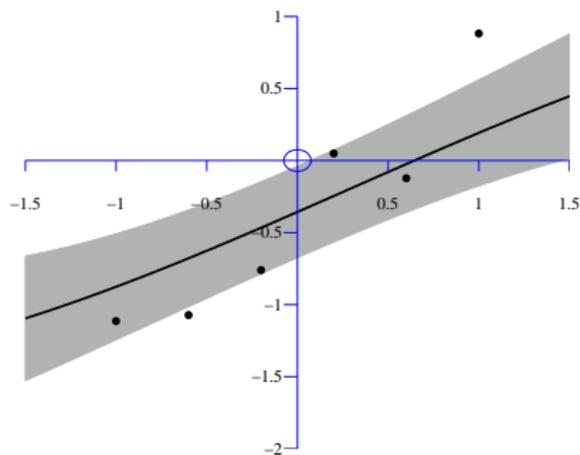
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

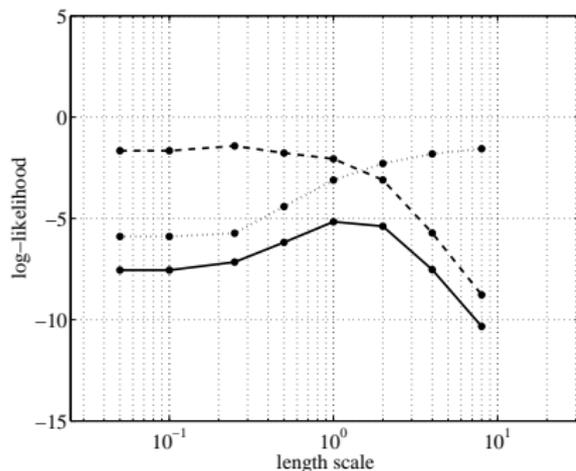
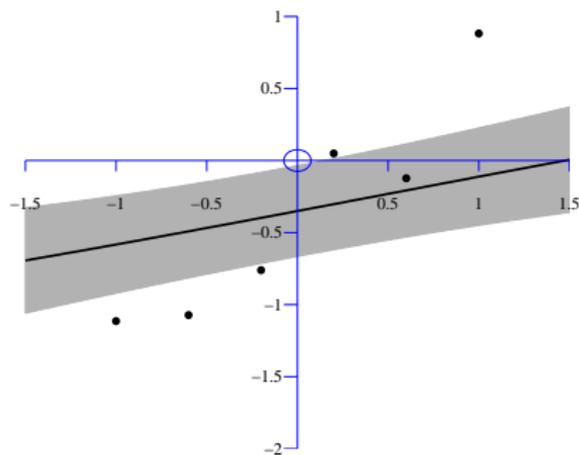
demOptimiseKern



Learning Kernel Parameters

Can we determine length scales and noise levels from the data?

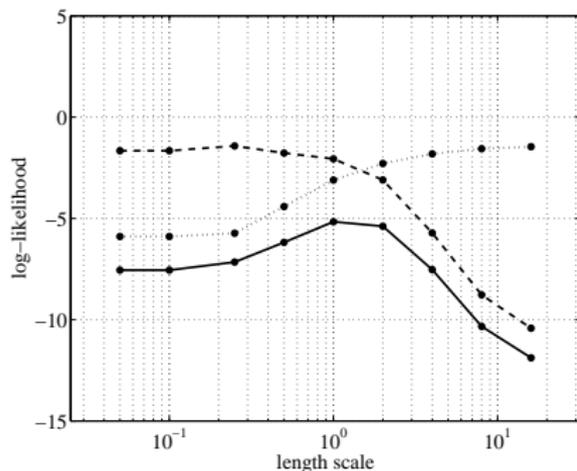
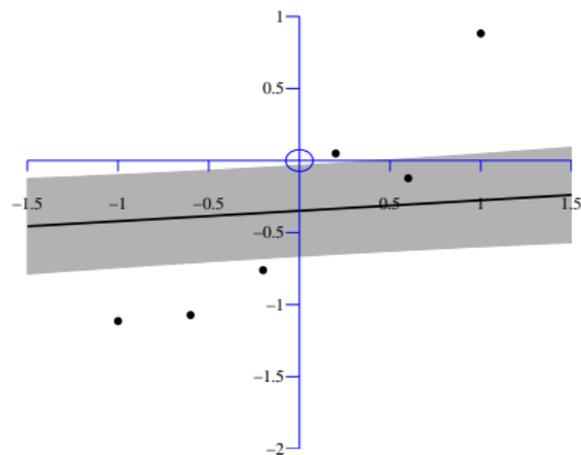
demOptimiseKern



Learning Kernel Parameters

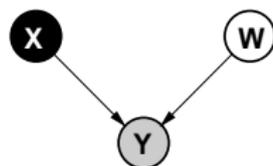
Can we determine length scales and noise levels from the data?

demOptimiseKern



Dual Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:
 - Define Gaussian prior over *parameters*, \mathbf{W} .
 - Integrate out *parameters*.



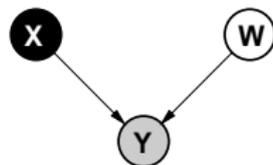
$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^n N(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{i=1}^D N(\mathbf{w}_{i,:} | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA

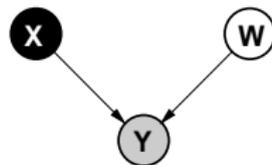
- Inspection of the marginal likelihood shows ...
 - The covariance matrix is a covariance function.



$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I})$$

Dual Probabilistic PCA

- Inspection of the marginal likelihood shows ...
 - The covariance matrix is a covariance function.
 - We recognise it as the 'linear kernel'.

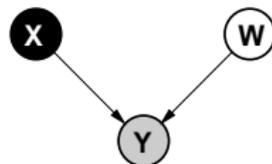


$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D \mathcal{N}(y_{:,j} | \mathbf{0}, \mathbf{K})$$

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

Dual Probabilistic PCA

- Inspection of the marginal likelihood shows ...
 - The covariance matrix is a covariance function.
 - We recognise it as the 'linear kernel'.



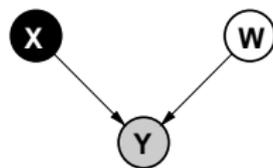
$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D N(y_{:,j}|\mathbf{0}, \mathbf{K})$$

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

This is a product of Gaussian processes with linear kernels.

Dual Probabilistic PCA

- Inspection of the marginal likelihood shows ...
 - The covariance matrix is a covariance function.
 - We recognise it as the 'linear kernel'.



$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^D N(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{K})$$

$$\mathbf{K} = ?$$

Replace linear kernel with non-linear kernel for non-linear model.

This is called the Gaussian Process Latent Variable Model (GPLVM)

RBF Kernel

- The RBF kernel has the form $k_{ij} = k(\mathbf{x}_{i,:}, \mathbf{x}_{j,:})$, where

$$k(\mathbf{x}_{i,:}, \mathbf{x}_{j,:}) = \alpha \exp\left(-\frac{(\mathbf{x}_{i,:} - \mathbf{x}_{j,:})^T (\mathbf{x}_{i,:} - \mathbf{x}_{j,:})}{2l^2}\right).$$

- No longer possible to optimise wrt \mathbf{X} via an eigenvalue problem.
- Instead find gradients with respect to \mathbf{X} , α , l and σ^2 and optimise using gradient methods.

'Swiss Roll'

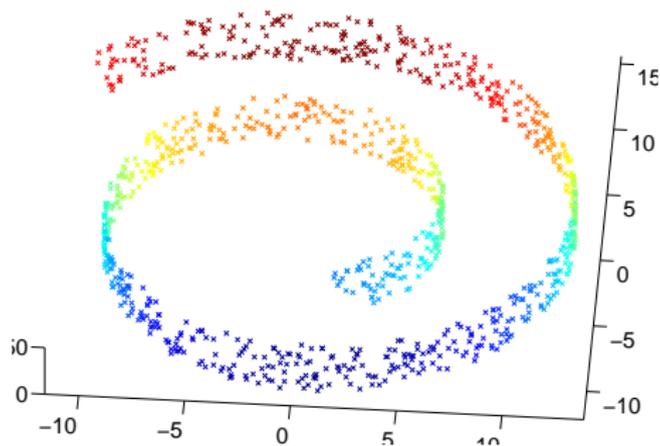


Figure: The 'Swiss Roll' data set is data in three dimensions that is inherently two dimensional.

Quality of solution is Initialisation Dependent

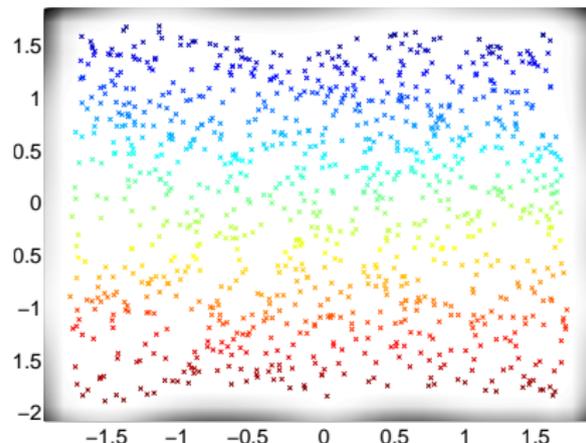
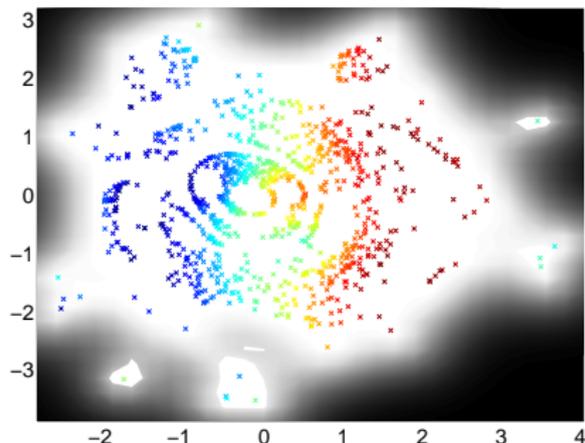


Figure: *Left:* Swiss roll solution initialised by PCA. *Right:* Swiss roll solution initialised by Isomap.

Stick Man Data

- $N = 55$ frames of motion capture.
- xyz locations of 34 points on the body.
- $D = 102$ dimensional data.
- “Run 1” available from http://accad.osu.edu/research/mocap/mocap_data.htm.

Changing



Angle



of Run



demStick1

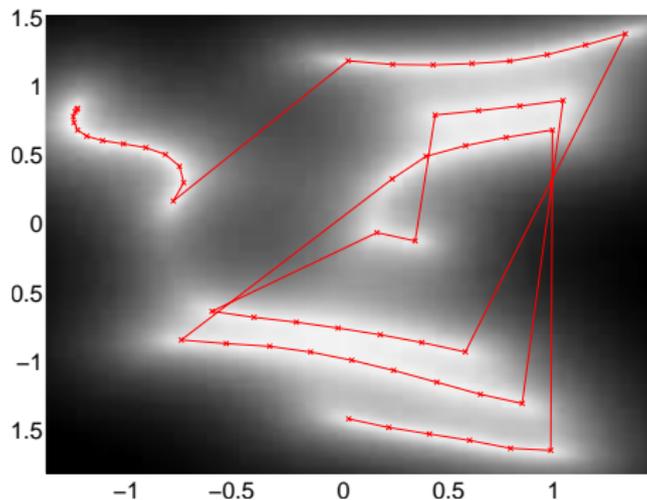


Figure: The latent space for the stick man motion capture data.

Non smooth latent spaces can be avoided by:

- Constrain the forward-mapping: using back-constraints
- Combine graph-based methods and non-linear latent variable models
- Use better optimization schemes that are less prone to get stuck in local minima
- Marginalize the latent space

Multi-Dimensional Scaling with a Mapping

- Lowe and Tipping (1997) made latent positions a function of the data.

$$x_{ij} = f_j(\mathbf{y}_i; \mathbf{w})$$

- Function was either multi-layer perceptron or a radial basis function network.
- Their motivation was different from ours:
 - They wanted to add the advantages of a true mapping to multi-dimensional scaling.

Back Constraints

- We can use the same idea to force the GP-LVM to respect local distances (Lawrence and Quinonero Candela, 2006).
- By constraining each \mathbf{x}_i to be a 'smooth' mapping from \mathbf{y}_i local distances can be respected.
- This works because in the GP-LVM we maximise wrt latent variables, we don't integrate out.
- Can use any 'smooth' function:
 - 1 Neural network.
 - 2 RBF Network.
 - 3 Kernel based mapping.

Computing Gradients

- GP-LVM normally proceeds by optimising

$$L(\mathbf{X}) = \log p(\mathbf{Y}|\mathbf{X})$$

with respect to \mathbf{X} using $\frac{dL}{d\mathbf{X}}$.

- The back constraints are of the form

$$x_{ij} = f_j(\mathbf{y}_{i,:}; \mathbf{B})$$

where \mathbf{B} are parameters.

- We can compute $\frac{dL}{d\mathbf{B}}$ via chain rule and optimise parameters of mapping.

Motion Capture Results

demStick1 and demStick3

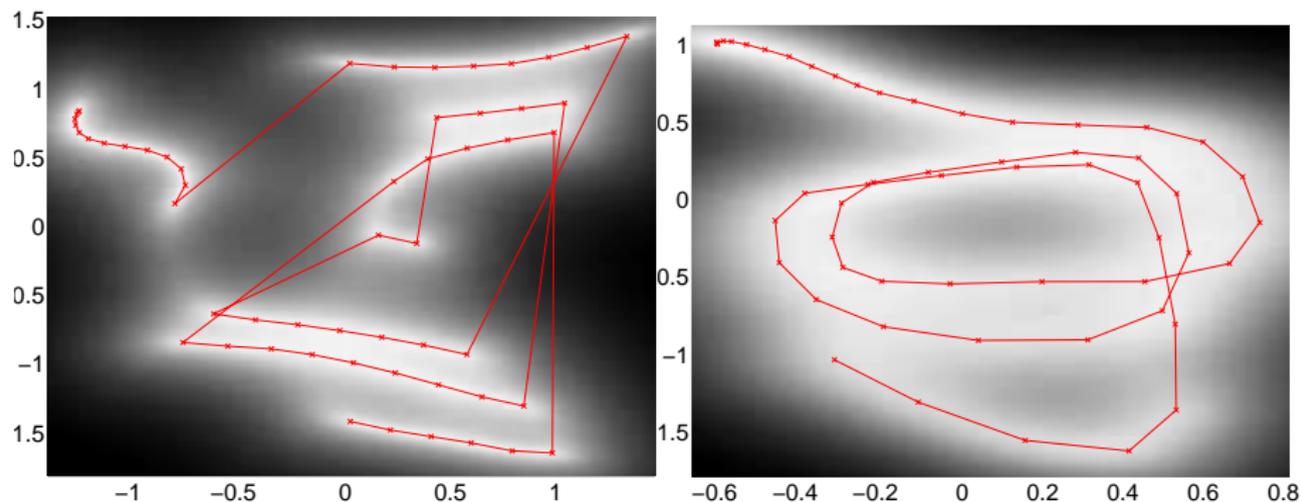
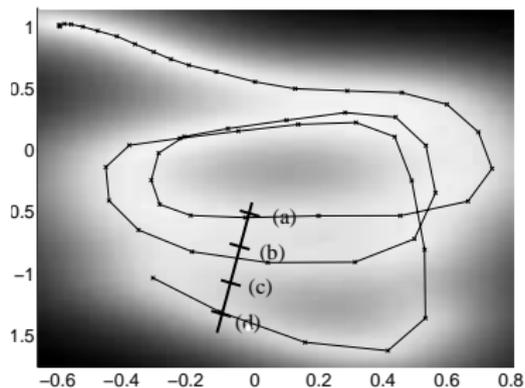


Figure: The latent space for the motion capture data with (*right*) and without (*left*) dynamics. The dynamics use a Gaussian process with an RBF kernel.

Stick Man Results

demStickResults



Projection into data space from four points in the latent space. The inclination of the runner changes becoming more upright.

Incorporating prior knowledge

- It is useful to use prior knowledge when additional information is available, e.g., cyclic motions, smoothness.
- We design priors over the latent space that incorporate the prior knowledge.
- Our prior is based on the Locally Linear Embedding (LLE) [Roweis, 01] cost function

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) + \lambda \sum_{i=1}^N \sum_{q=1}^d \left\| \mathbf{x}_{i,q} - \sum_{j \in \eta_i} w_{ij,q} \mathbf{x}_{j,q} \right\|^2$$

with $\mathbf{x}_{i,q}$ the q -th dimension of \mathbf{x}_i .

- We define the weights to reflect the prior knowledge.
- This is the Locally Linear GPLVM (LL-GPLVM) (Urtasun et al., 2008)

Incorporating prior knowledge

- It is useful to use prior knowledge when additional information is available, e.g., cyclic motions, smoothness.
- We design priors over the latent space that incorporate the prior knowledge.
- Our prior is based on the Locally Linear Embedding (LLE) [Roweis, 01] cost function

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) + \lambda \sum_{i=1}^N \sum_{q=1}^d \left\| \mathbf{x}_{i,q} - \sum_{j \in \eta_i} w_{ij,q} \mathbf{x}_{j,q} \right\|^2$$

with $\mathbf{x}_{i,q}$ the q -th dimension of \mathbf{x}_i .

- We define the weights to reflect the prior knowledge.
- This is the Locally Linear GPLVM (LL-GPLVM) (Urtasun et al., 2008)

Generate animations by sampling

- We learn style-content separation models using the following sources of prior knowledge (Urtasun et al. 2008)
 - ▶ smoothness: points close in observation space should be close in latent space.
 - ▶ cyclic structure: points with similar phase should be close.
 - ▶ transitions: points where a transition could happen should be close in the latent space.

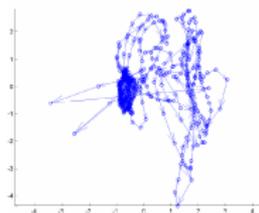


Figure: GPLVM

Figure: Topologies

Figure: Sampling

Problems with the GPLVM

- It relies on the optimization of a non-convex function

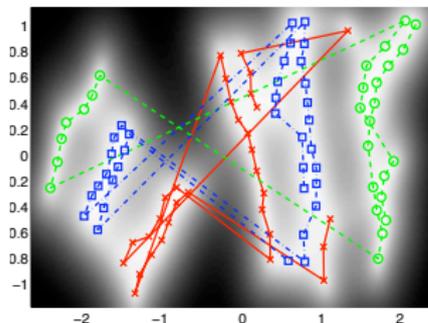
$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) .$$

Problems with the GPLVM

- It relies on the optimization of a non-convex function

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) .$$

- Even with the right dimensionality, they can result in poor representations if initialized far from the optimum.

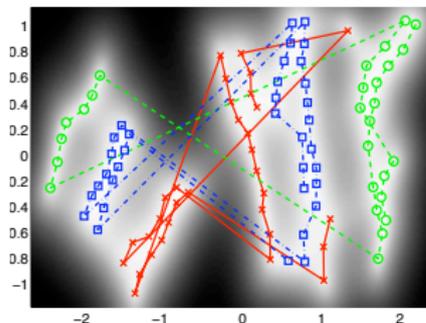


Problems with the GPLVM

- It relies on the optimization of a non-convex function

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) .$$

- Even with the right dimensionality, they can result in poor representations if initialized far from the optimum.



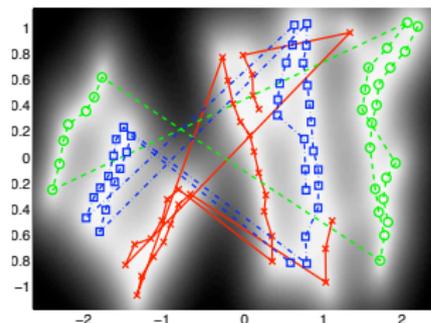
- This is even worse if the dimensionality of the latent space is small.

Problems with the GPLVM

- It relies on the optimization of a non-convex function

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T).$$

- Even with the right dimensionality, they can result in poor representations if initialized far from the optimum.



- This is even worse if the dimensionality of the latent space is small.
- As a consequence this models have only been applied to small databases of a single activity.

- No distortion is introduced by an initialization step; the latent coordinates are initialized to be the original observations

$$\mathbf{X}_{init} = \mathbf{Y}$$

- We introduce a prior over the latent space that encourages latent spaces to be low dimensional.
- Our method is able to estimate the latent space and its dimensionality (Geiger et al., 2009).

Continuous dimensionality reduction

- We want to encourage latent space that are low-dimensional.
- Dimensionality can be measure by the rank of $\mathbf{X}\mathbf{X}^T$.

Continuous dimensionality reduction

- We want to encourage latent space that are low-dimensional.
- Dimensionality can be measure by the rank of \mathbf{XX}^T .
- We would like to penalize the rank, but the rank is a discrete function. The optimization would have to solve a complex combinatorial problem.

Continuous dimensionality reduction

- We want to encourage latent space that are low-dimensional.
- Dimensionality can be measure by the rank of $\mathbf{X}\mathbf{X}^T$.
- We would like to penalize the rank, but the rank is a discrete function. The optimization would have to solve a complex combinatorial problem.
- We relax the rank minimization and define a prior that encourages sparsity of the eigenvalues, such that:

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^T) + \alpha \sum_{i=1}^D \phi(s_i)$$

with s_i the eigenvalues of $\bar{\mathbf{X}}\bar{\mathbf{X}}^T$, $\bar{\mathbf{X}}$ the zero-mean \mathbf{X} , and ϕ is a function that encourages sparsity.

Continuous dimensionality reduction

- We want to encourage latent space that are low-dimensional.
- Dimensionality can be measure by the rank of $\mathbf{X}\mathbf{X}^T$.
- We would like to penalize the rank, but the rank is a discrete function. The optimization would have to solve a complex combinatorial problem.
- We relax the rank minimization and define a prior that encourages sparsity of the eigenvalues, such that:

$$\mathcal{L} = \frac{D}{2} \ln |\mathbf{K}| + \frac{D}{2} \text{tr}(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^T) + \alpha \sum_{i=1}^D \phi(s_i)$$

with s_i the eigenvalues of $\bar{\mathbf{X}}\bar{\mathbf{X}}^T$, $\bar{\mathbf{X}}$ the zero-mean \mathbf{X} , and ϕ is a function that encourages sparsity.

Choice of the penalty function

- Common choice for sparseness is the power family

$$\phi(s_i, \rho) = |s_i|^\rho$$

$\rho = 1$ is a Laplace prior (i.e., L1 norm), which is linear.

- However, our objective function is non-convex. We use a penalty that drives faster to zero the small singular values

$$\phi(s_i) = \log(1 + \beta s_i) .$$

Choice of the penalty function

- Common choice for sparseness is the power family

$$\phi(s_i, \rho) = |s_i|^\rho$$

$\rho = 1$ is a Laplace prior (i.e., L1 norm), which is linear.

- However, our objective function is non-convex. We use a penalty that drives faster to zero the small singular values

$$\phi(s_i) = \log(1 + \beta s_i) .$$

Estimating the dimensionality

- Minimizing the negative log posterior results in a reduction of the energy of the spectrum. We prevent this by optimizing instead

$$\min \mathcal{L} \quad \text{s.t.} \quad \forall i \ s_i \geq 0, \quad E(\mathbf{Y}) - E(\mathbf{X}) = 0,$$

with the energy $E(\mathbf{X}) = \sum_i s_i^2$.

- Finally, we choose the dimensionality to be

$$Q = \operatorname{argmax}_i \frac{s_i}{s_{i+1} + \epsilon}$$

where $\epsilon \ll 1$, and $s_1 \geq s_2 \geq \dots \geq s_D$

Estimating the dimensionality

- Minimizing the negative log posterior results in a reduction of the energy of the spectrum. We prevent this by optimizing instead

$$\min \mathcal{L} \quad \text{s.t.} \quad \forall i \ s_i \geq 0, \quad E(\mathbf{Y}) - E(\mathbf{X}) = 0,$$

with the energy $E(\mathbf{X}) = \sum_i s_i^2$.

- Finally, we choose the dimensionality to be

$$Q = \operatorname{argmax}_i \frac{s_i}{s_{i+1} + \epsilon}$$

where $\epsilon \ll 1$, and $s_1 \geq s_2 \geq \dots \geq s_D$

Dimensionality Estimation Results

Results on mocap

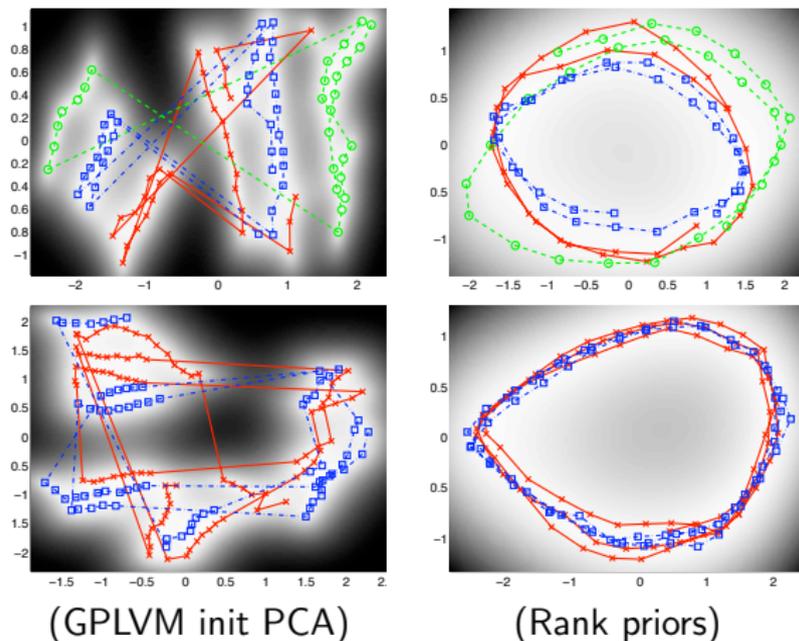


Figure: Running (top) and walking (bottom) models from mocap data. Different subjects are depicted in different colors. Unlike with the GPLVM, the latent coordinates using rank priors are very smooth.

Stacking Gaussian Processes

- Regressive dynamics provides a simple hierarchy.
- The input space of the GP is governed by another GP.
- By stacking GPs we can consider more complex hierarchies.
- Ideally we should marginalise latent spaces
 - In practice we seek MAP solutions.

Two Correlated Subjects

demHighFive1

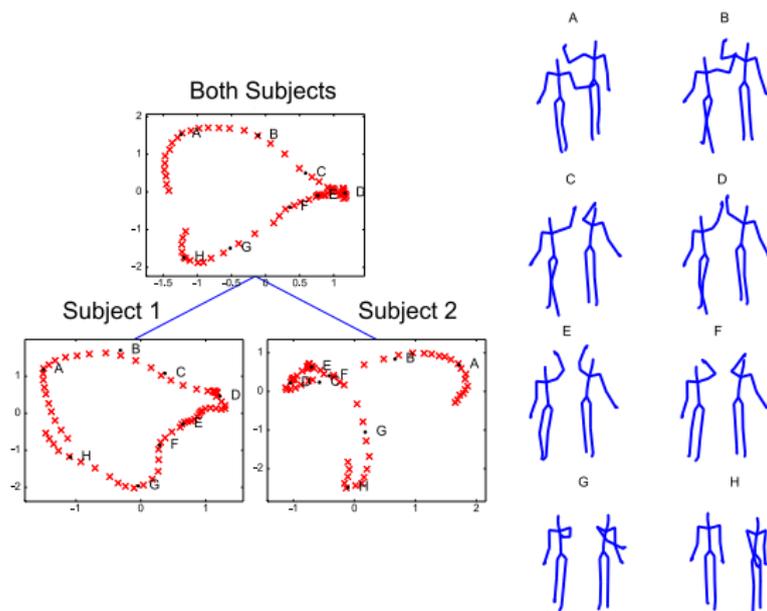


Figure: Hierarchical model of a 'high five'.

Decomposition of Body

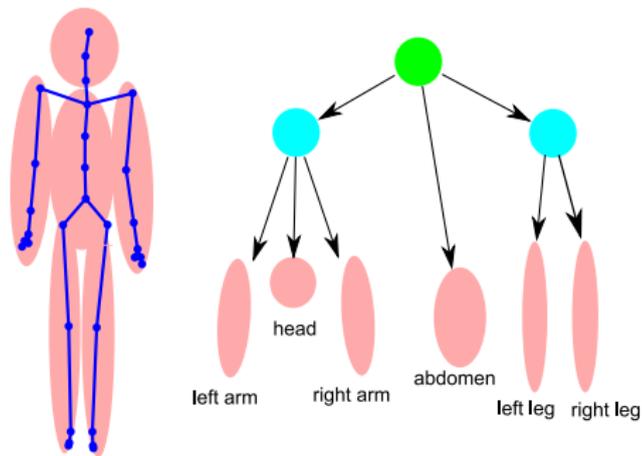


Figure: Decomposition of a subject.

Single Subject Run/Walk

demRunWalk1

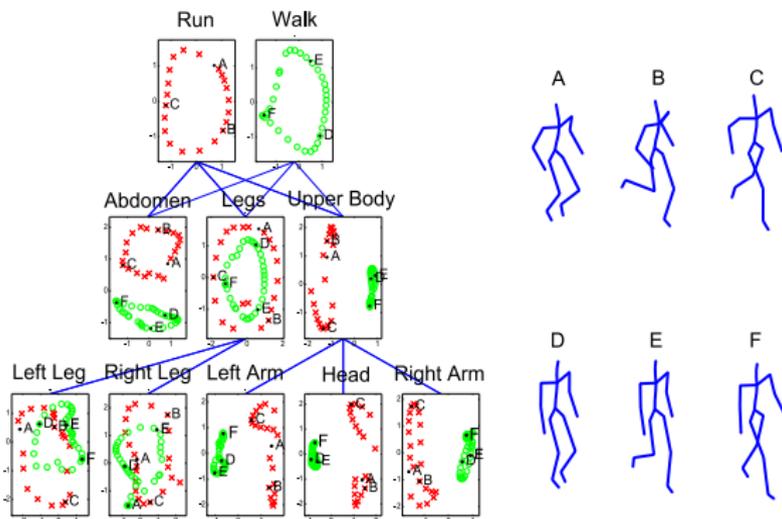


Figure: Hierarchical model of a walk and a run.

More?

- If you want to learn more, look at the additional material.
- Otherwise, do the research project on this topic!
- Next week we will do dynamical models.
- Let's do some exercises now!