



CS598PS – Machine Learning for Signal Processing

# Dimensionality Reduction – Non-linear Approaches

20 September 2017

# Today's lecture

- Non-linear dimensionality reduction
- Kernel PCA
- Manifold Methods

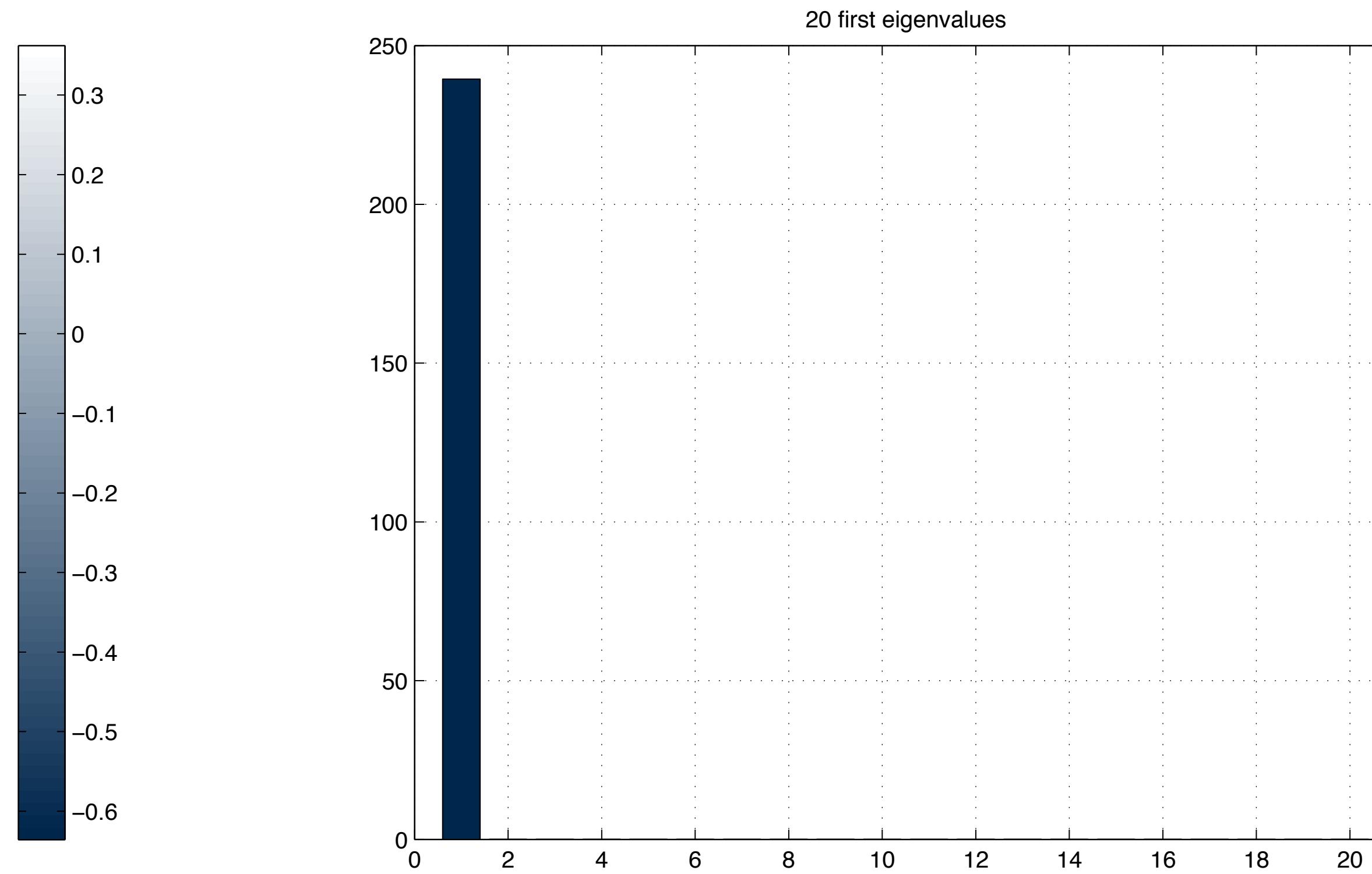
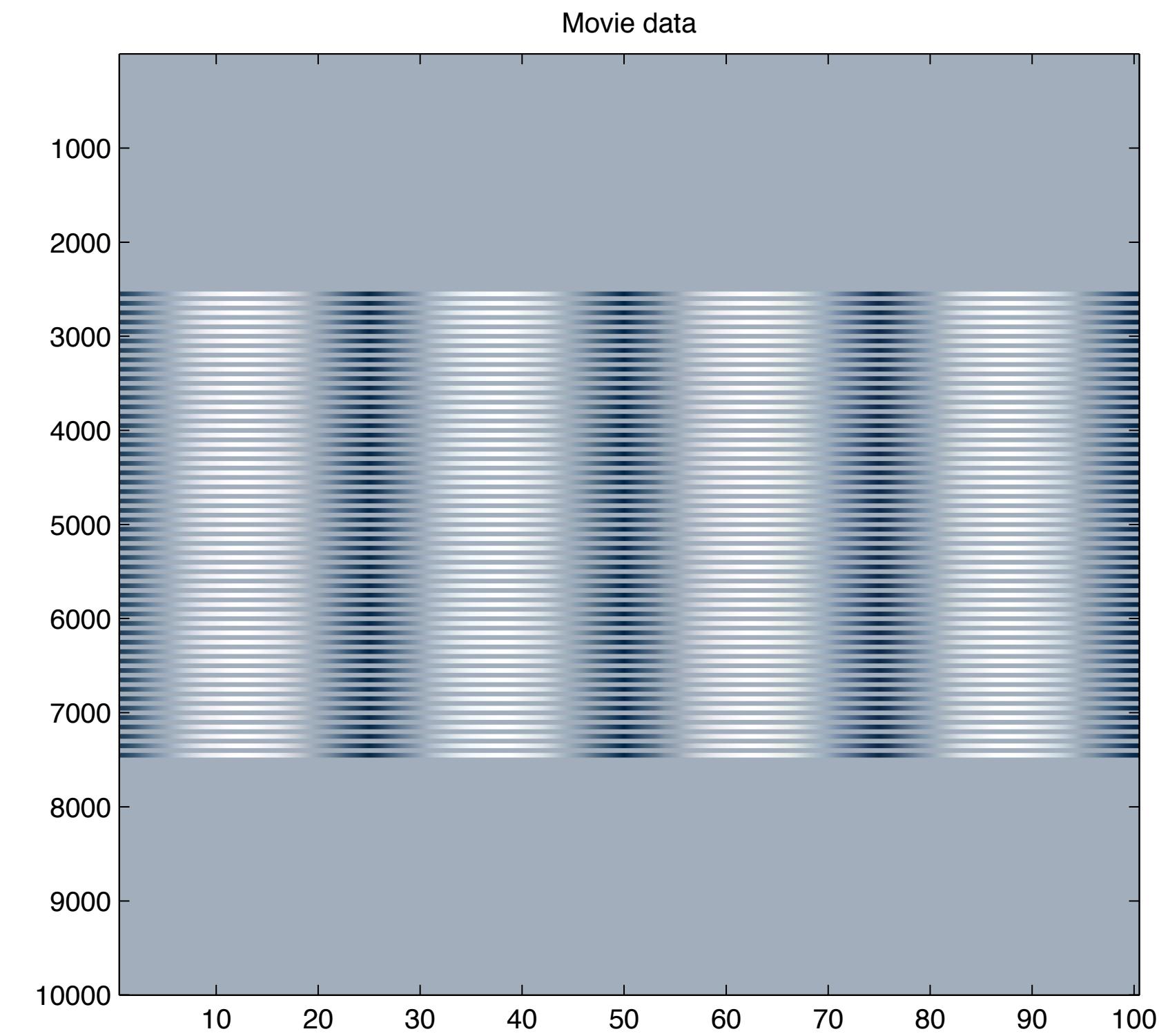
# Why dimensionality reduction

- Data can be hugely redundant
  - What are the (data) dimensions of this  $100 \times 100 \times 100$  video?



# Very low dimensionality!

- We effectively need only one pixel to express the video
  - All other active pixels behave the same



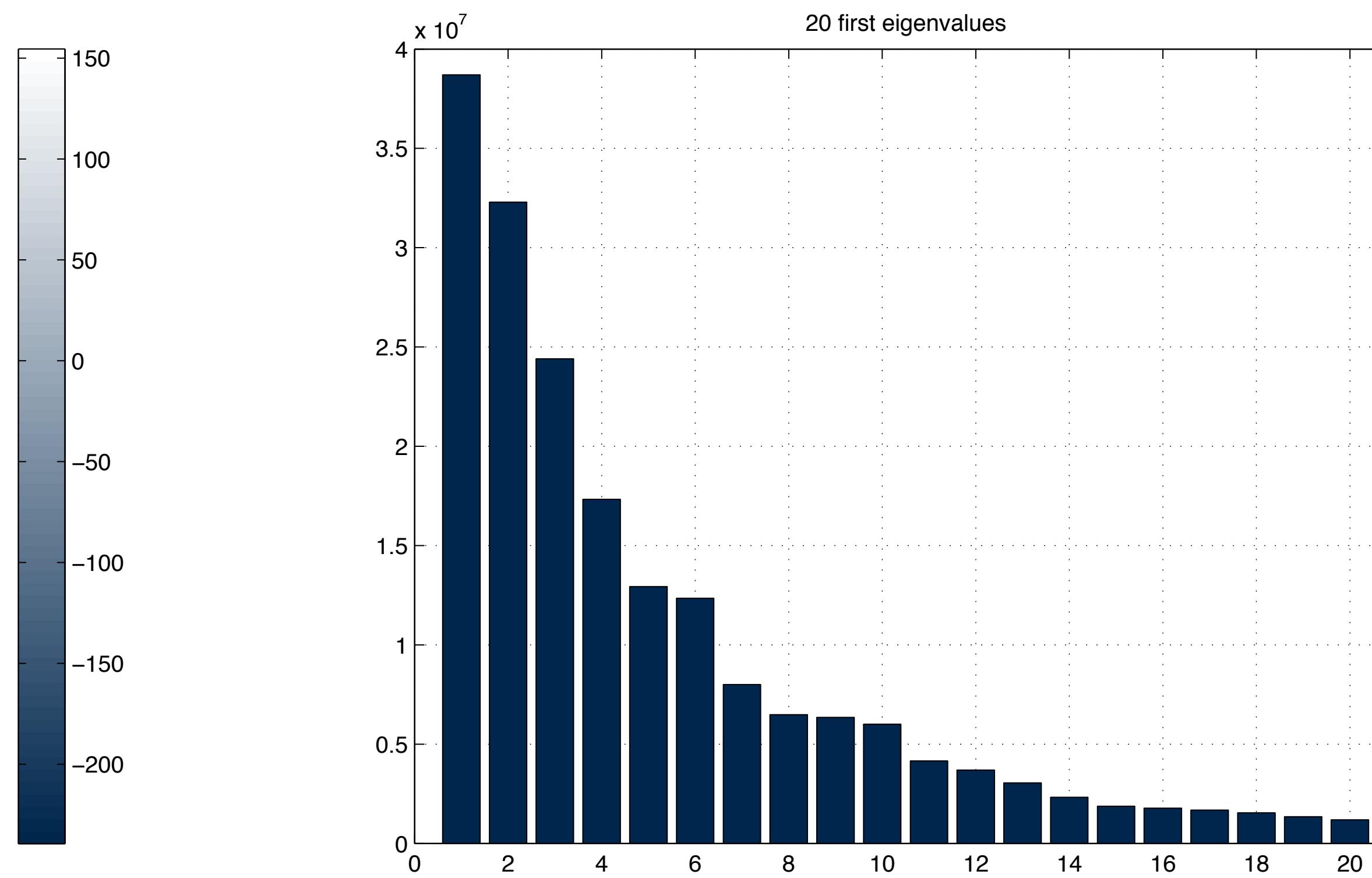
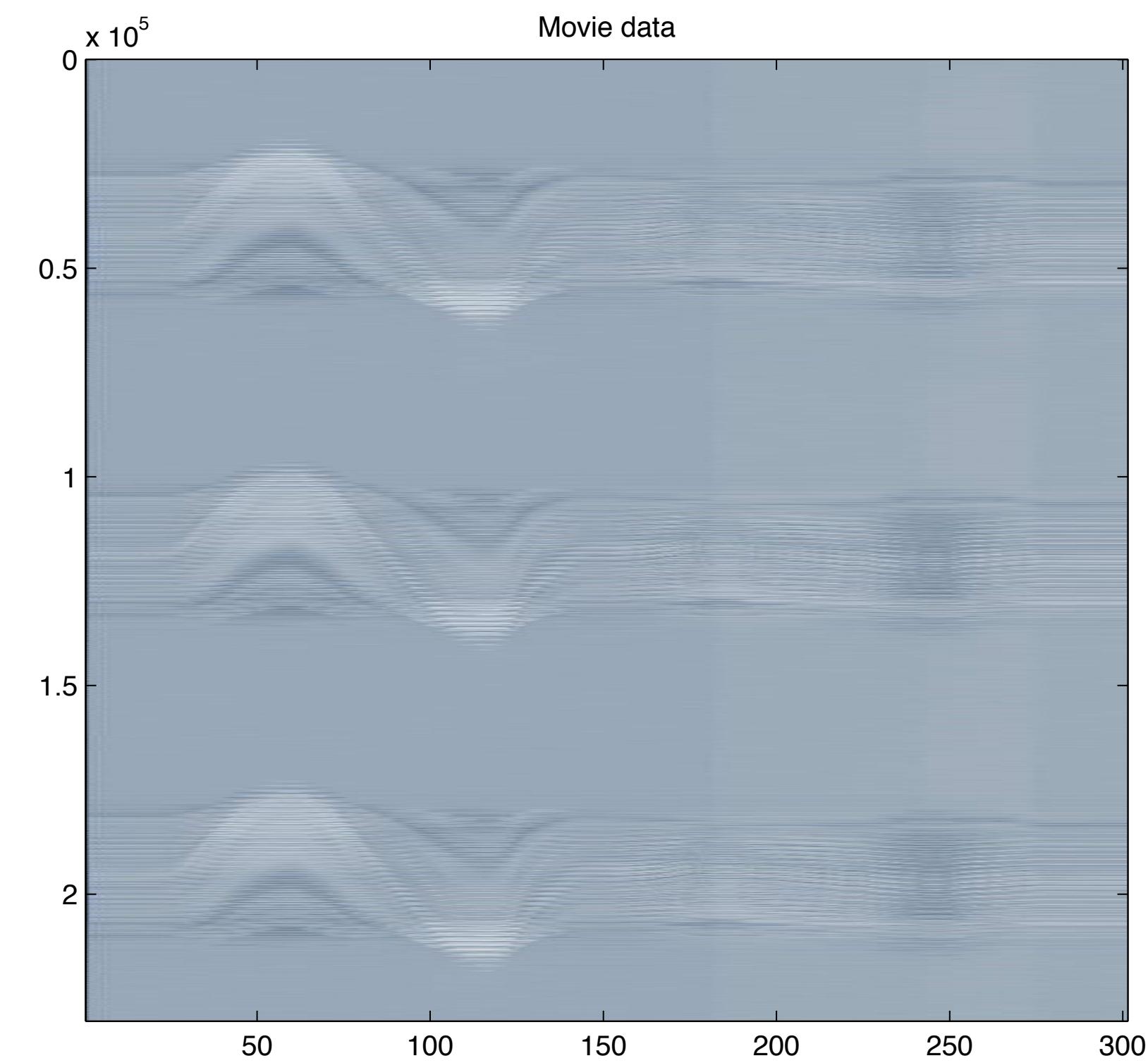
# In real life though ...

- What is the dimensionality now?
  - Lots of data,  $320 \times 240 \times 3 \times 337$



# It looks high dimensional

- Less than the data implies, but more than what we see
  - I still think it is around 4-dimensional, though

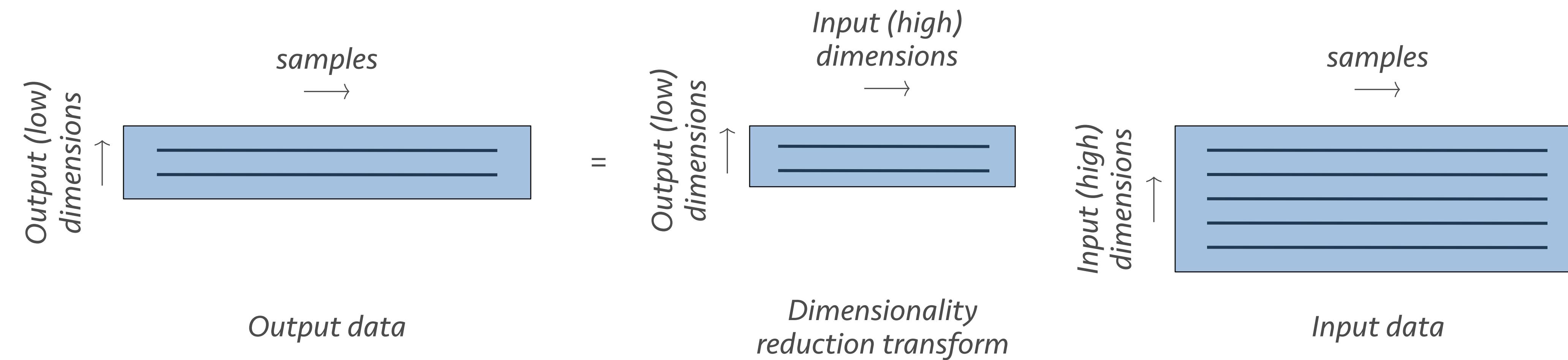


# Linear dimensionality reduction

- Linear transform to drop dimensions

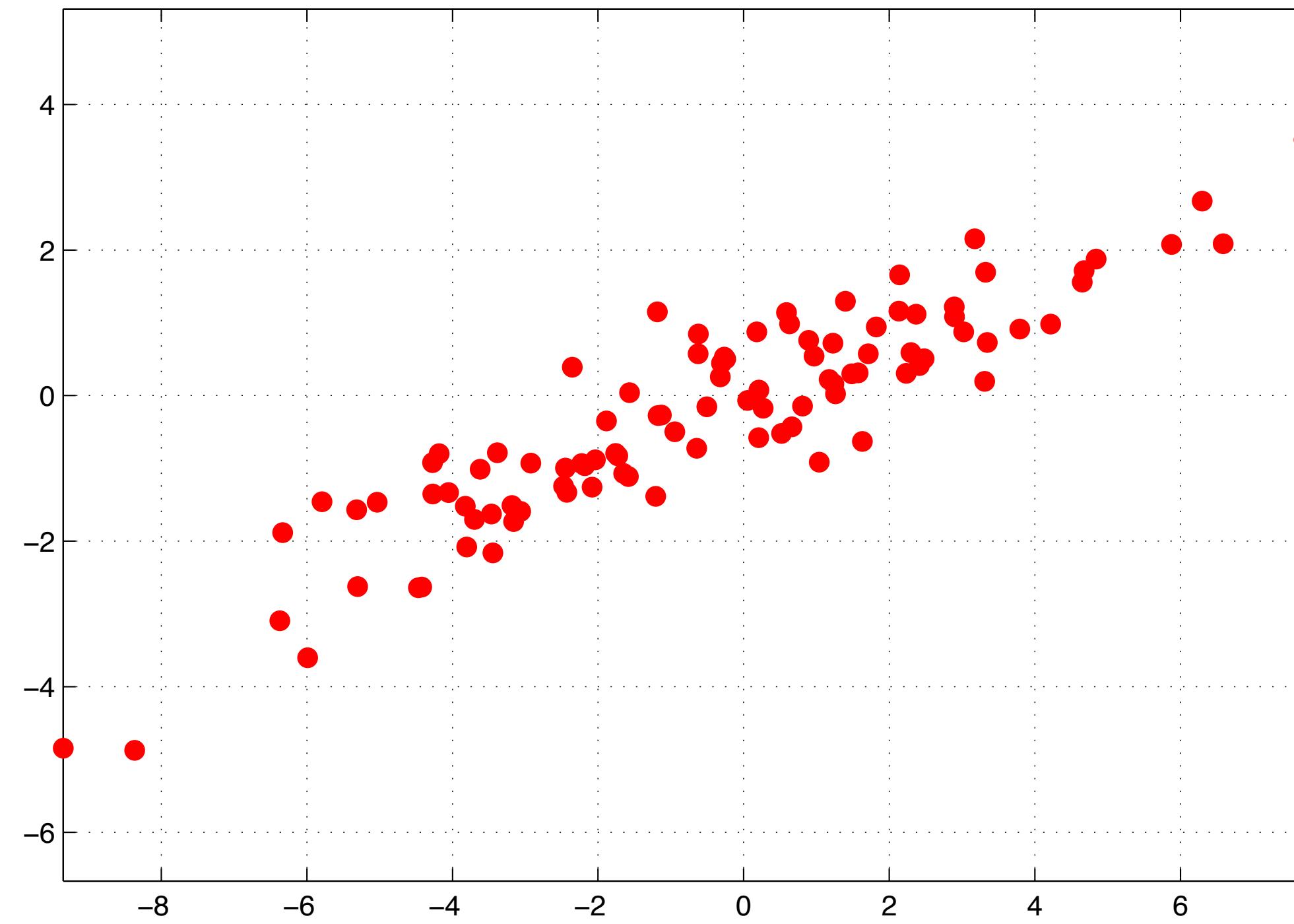
$$\mathbf{Y} = \mathbf{W} \cdot \mathbf{X}$$

- We saw this with PCA/NMF



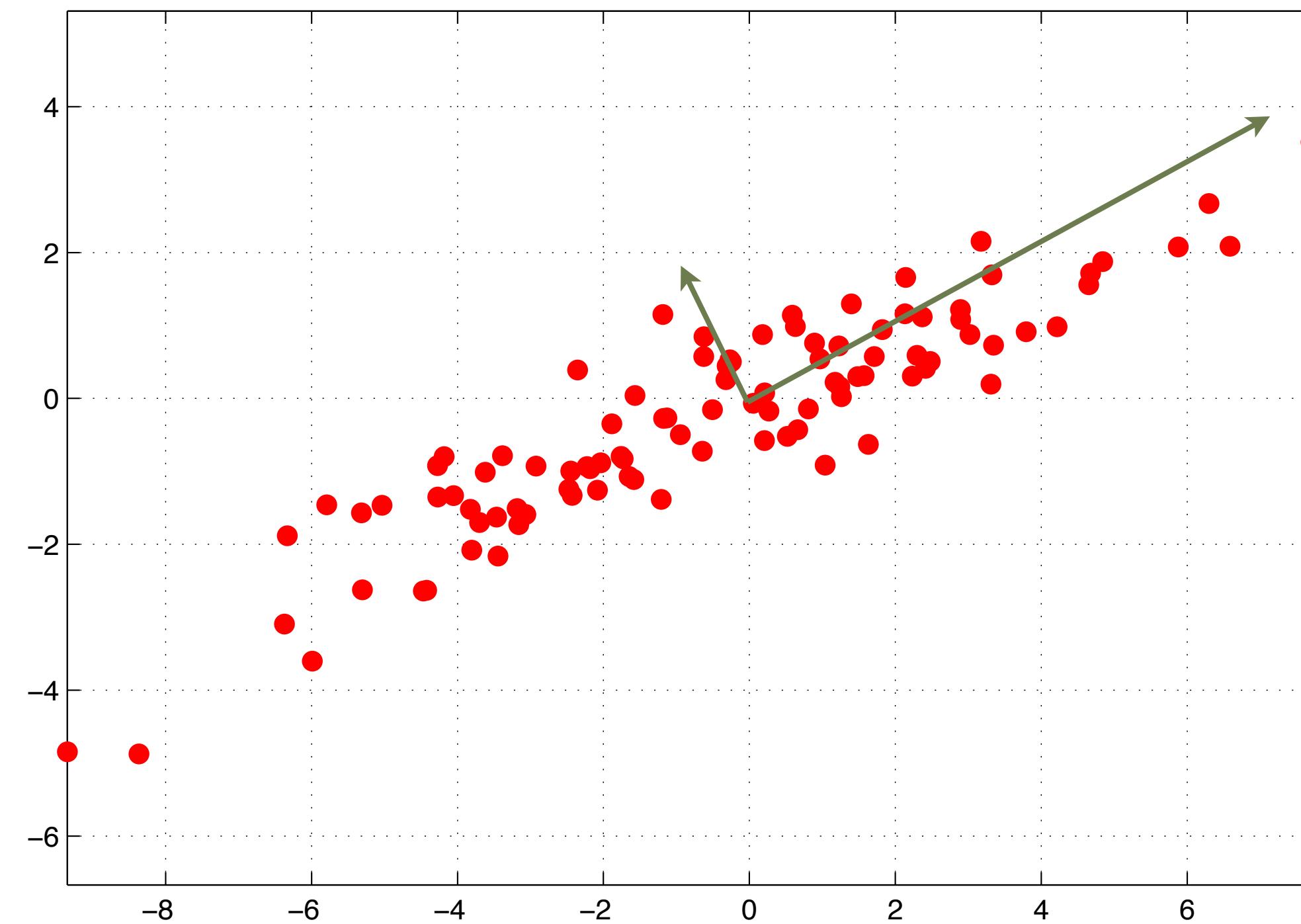
# “Easy” data

- Linear transforms work fine with simple data



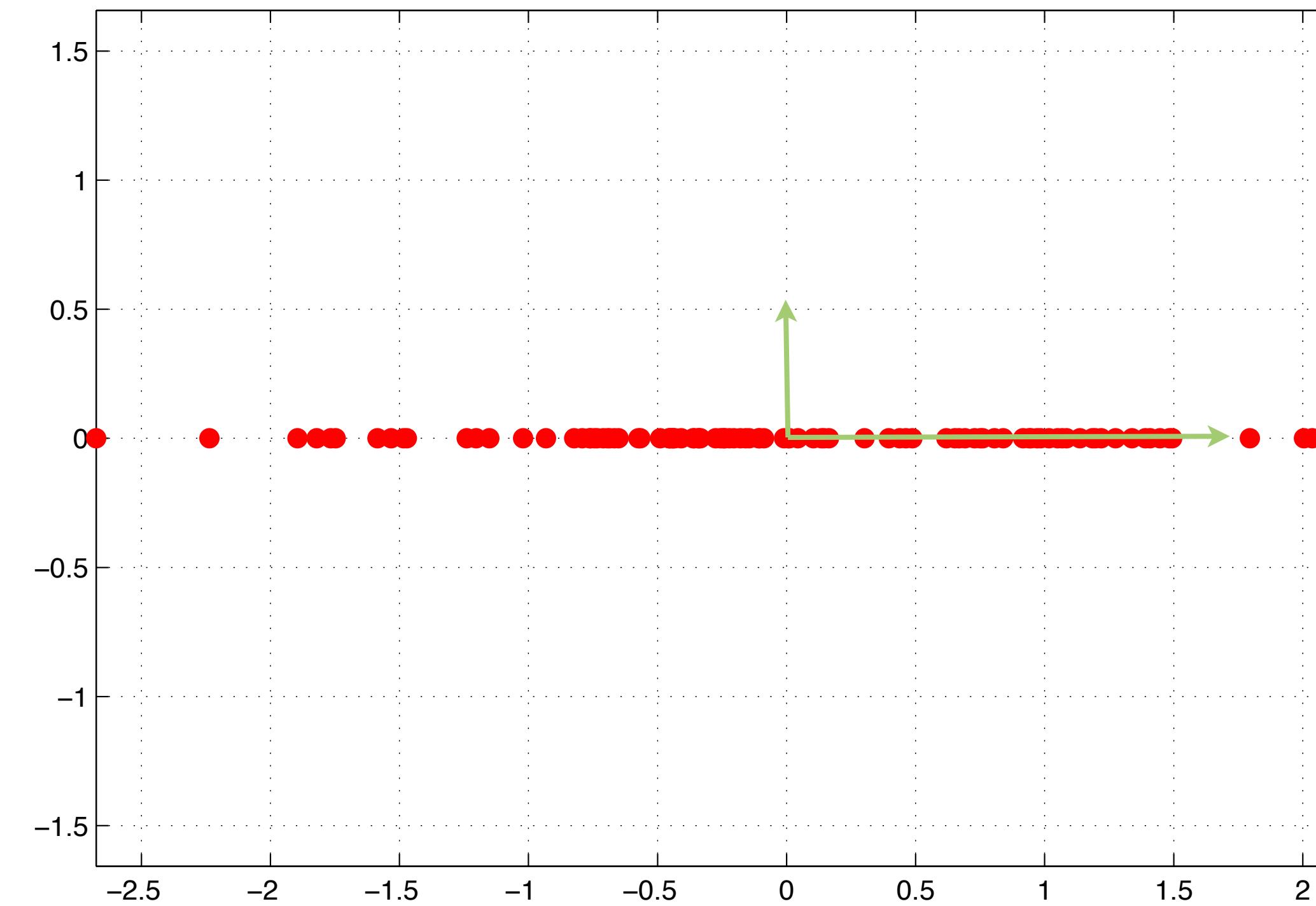
# “Easy” data

- Principal components show directions of maximal variance



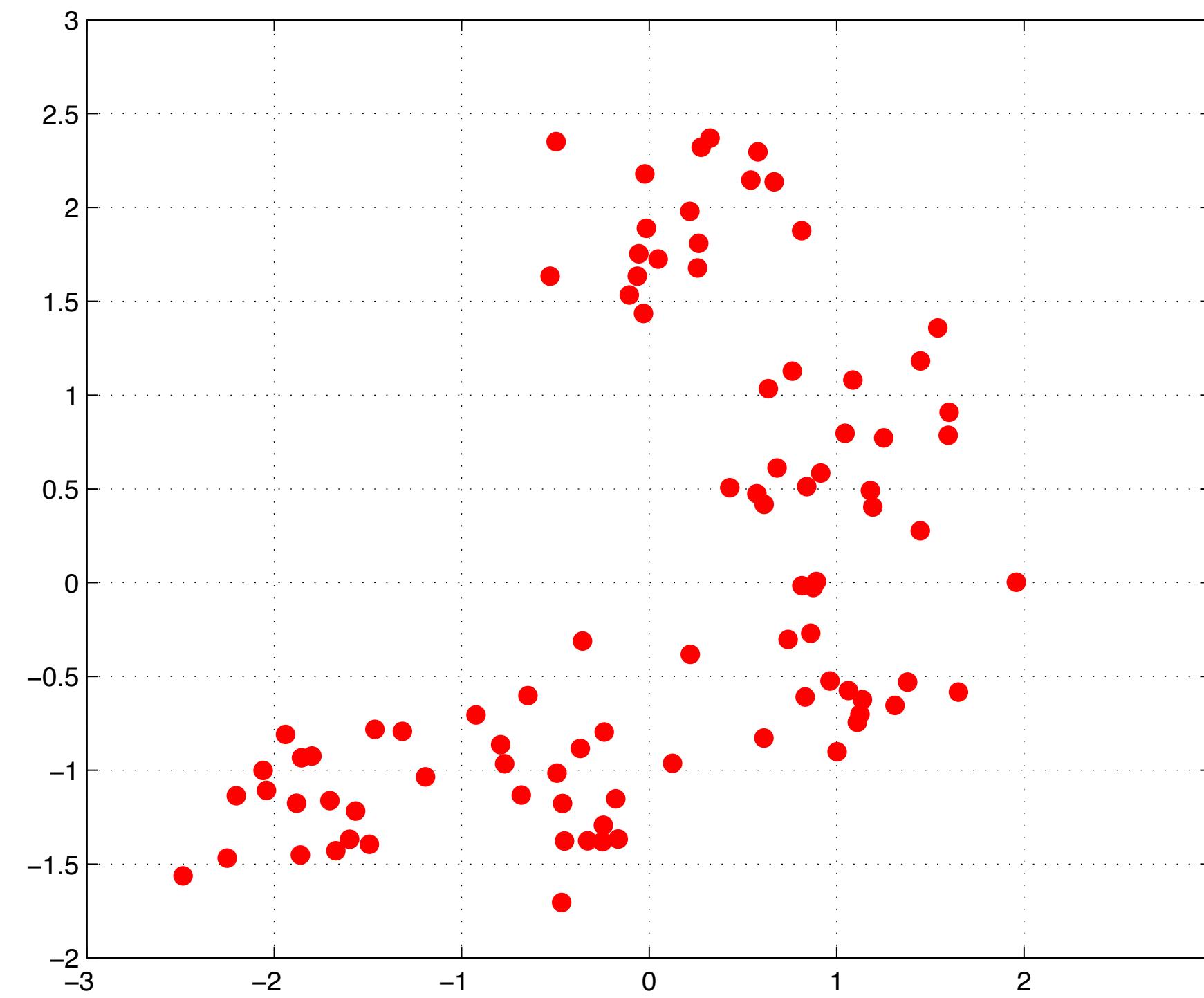
# “Easy” data

- Keep only the large dimensions to reduce dimensionality



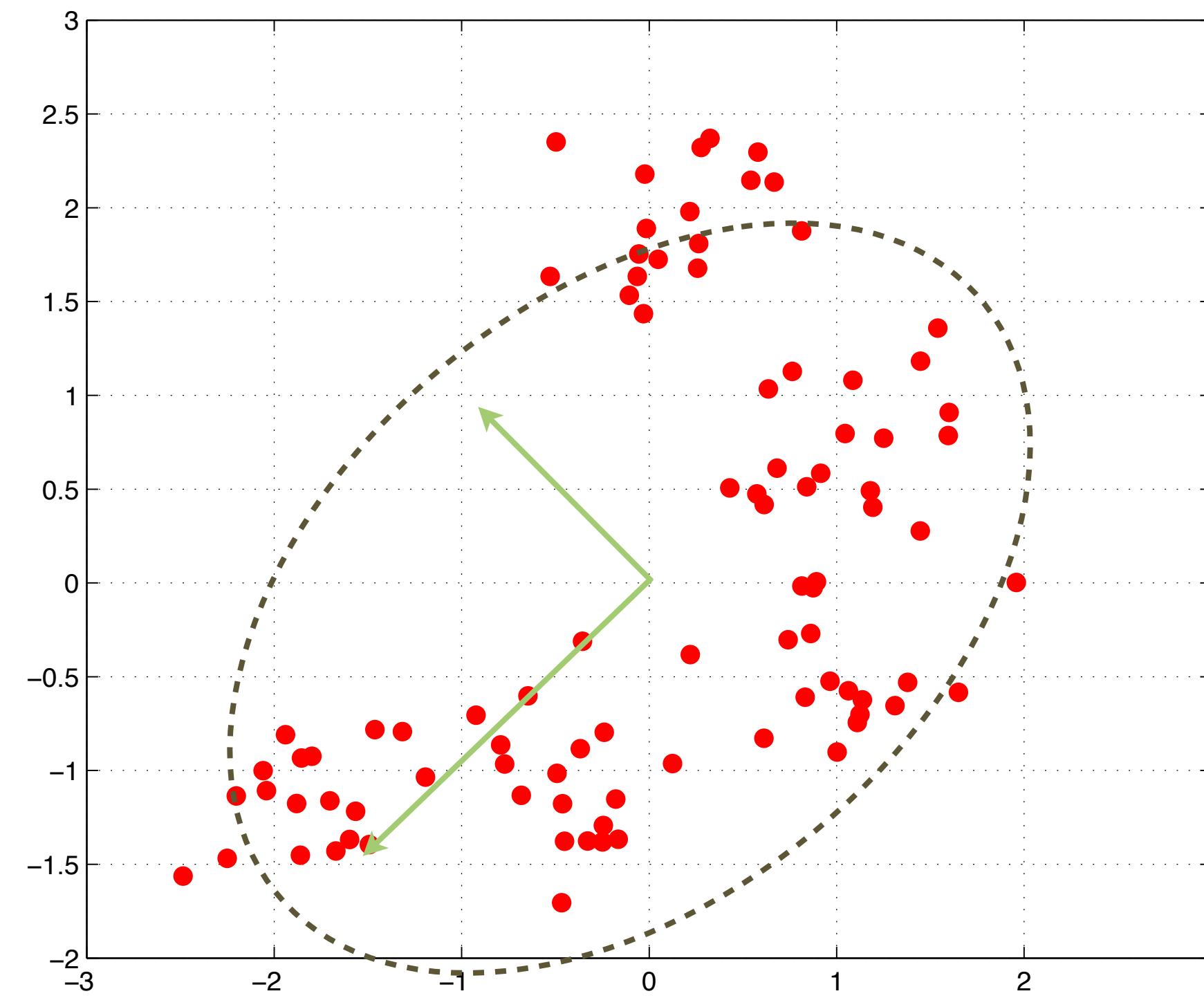
# More complicated data

- If data isn't Gaussian(ish) PCA doesn't help



# More complicated data

- Not much to interpret here ...

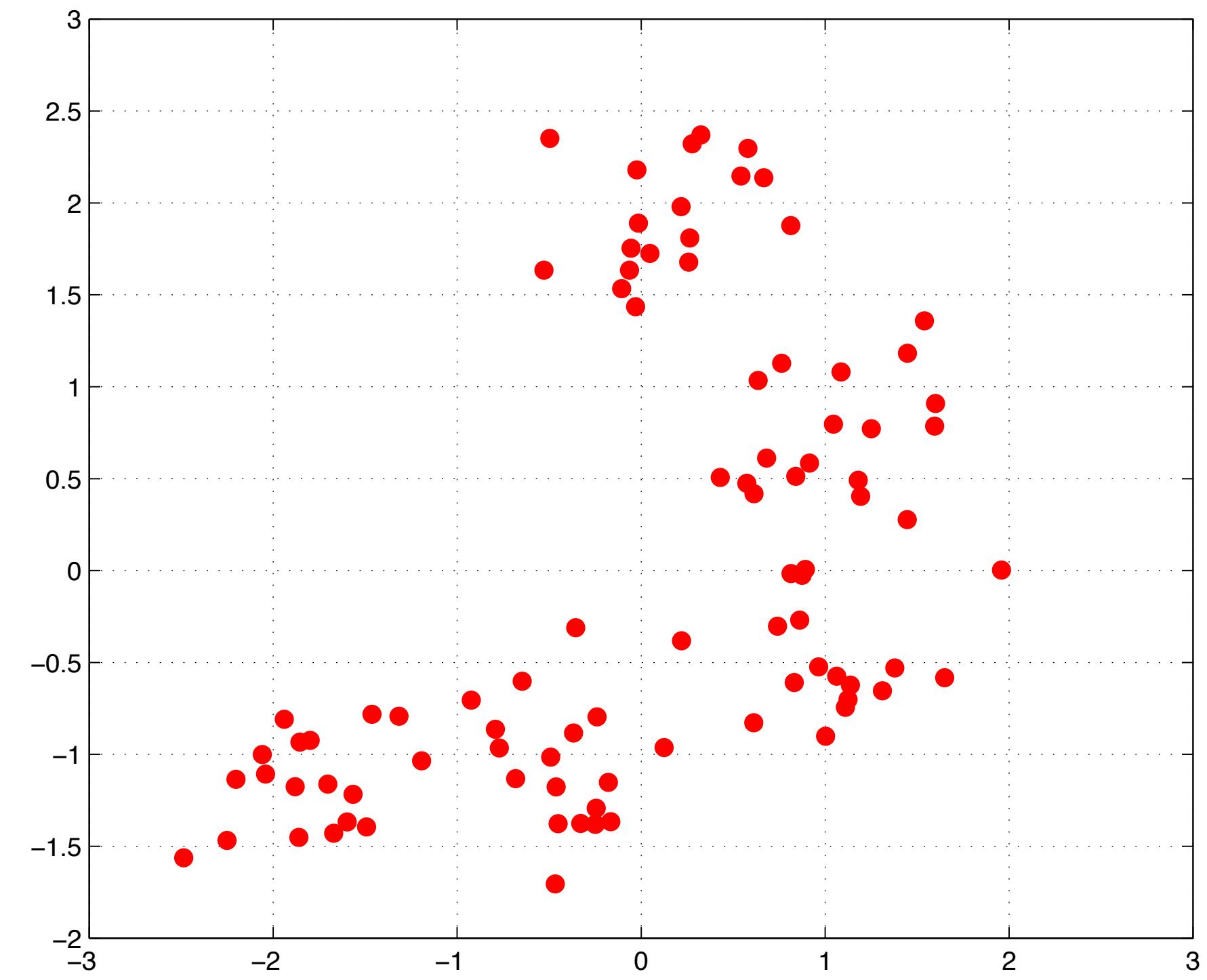


# Pinpointing the problem

- The principal components are linear
- The data we observe will not always conform to that
- Can we define “non-linear” components?

# Going the other way

- Suppose that we have some “curvy” data
- There should be a non-linear mapping that “straightens” that data out



# Going the wrong way

- With  $N, D$ -dimensional data points

$$\mathbf{x}_n \in \mathbb{R}^D, n = \{1, \dots, N\}$$

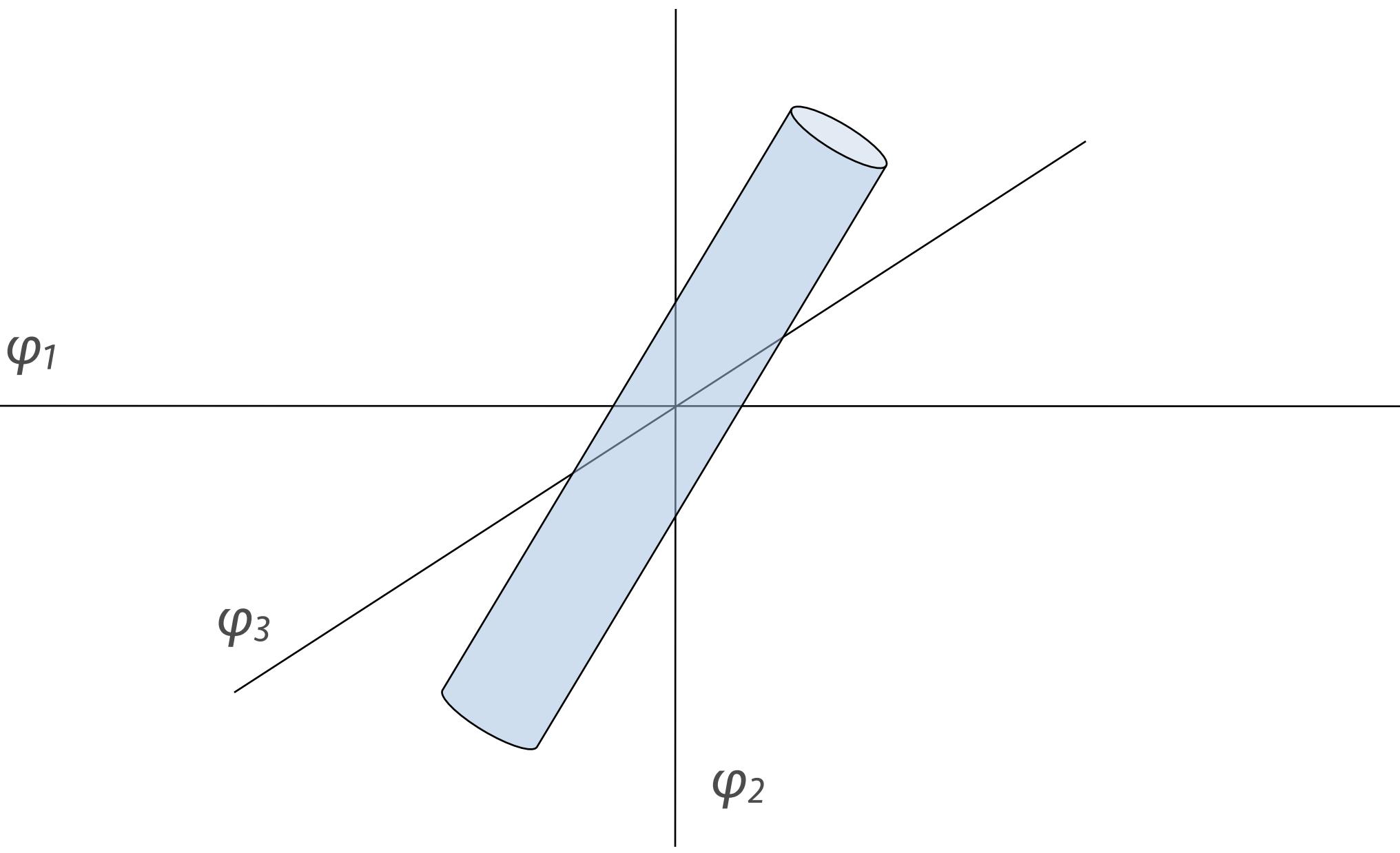
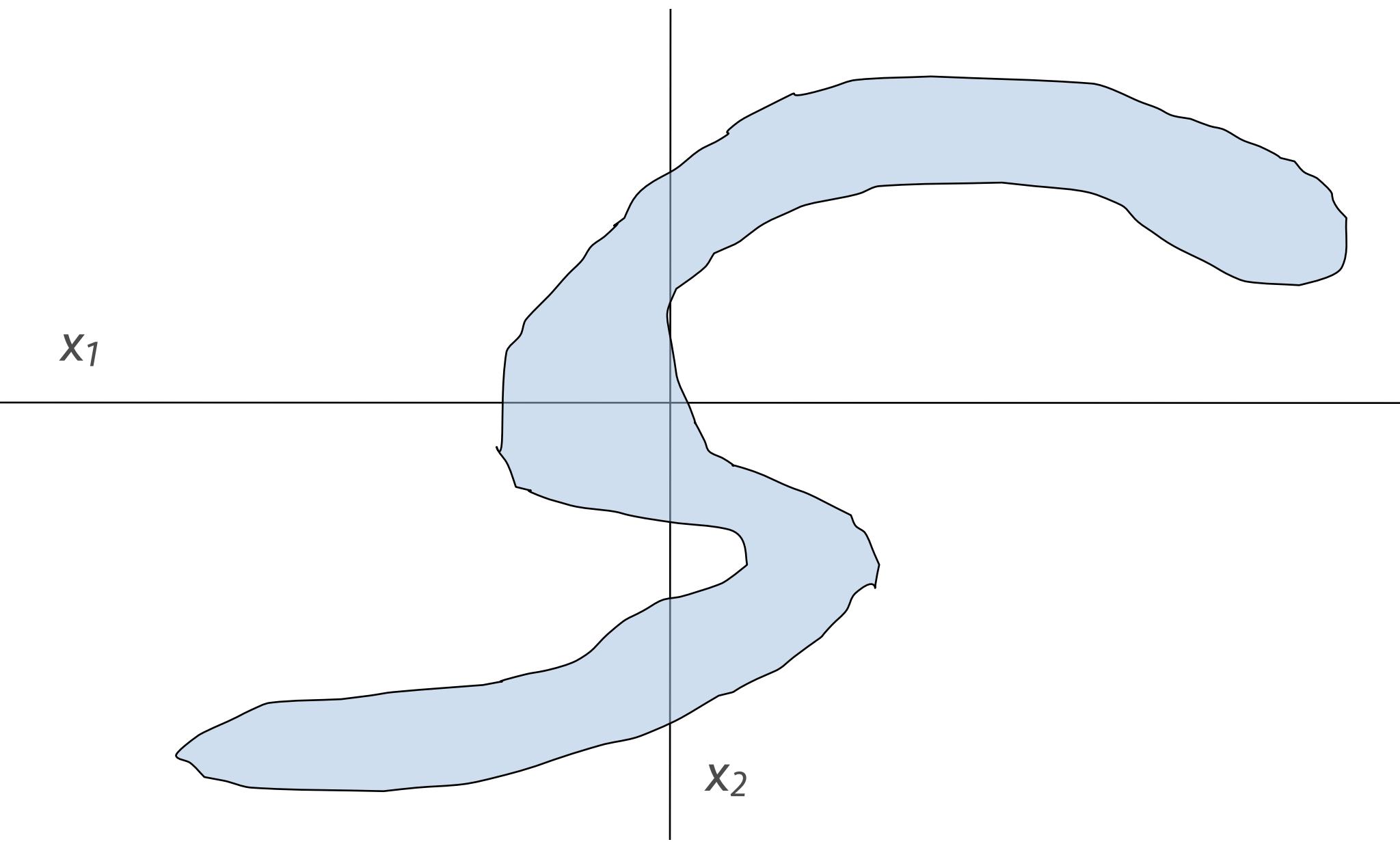
- Assume an unspecified non-linear mapping

$$\phi: \mathbb{R}^D \rightarrow \mathbb{R}^M, \mathbf{x} \mapsto \mathbf{z} = \phi(\mathbf{x}), M > D$$

- i.e. we non-linearize and increase the dimension of our data!

# In pictures

$$\mathbf{x} \in \mathbb{R}^2 \rightarrow \phi(\mathbf{x}) \in \mathbb{R}^3$$

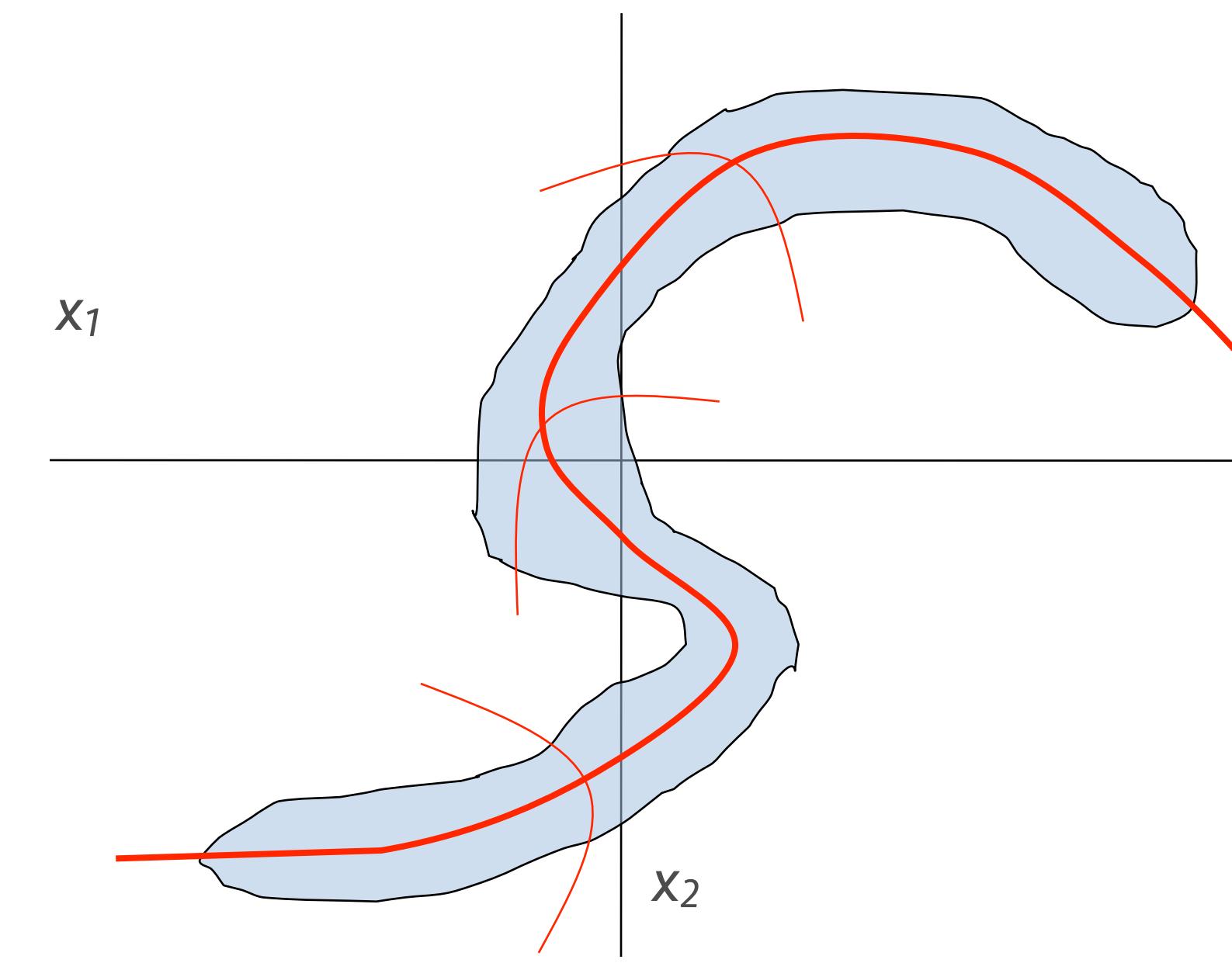
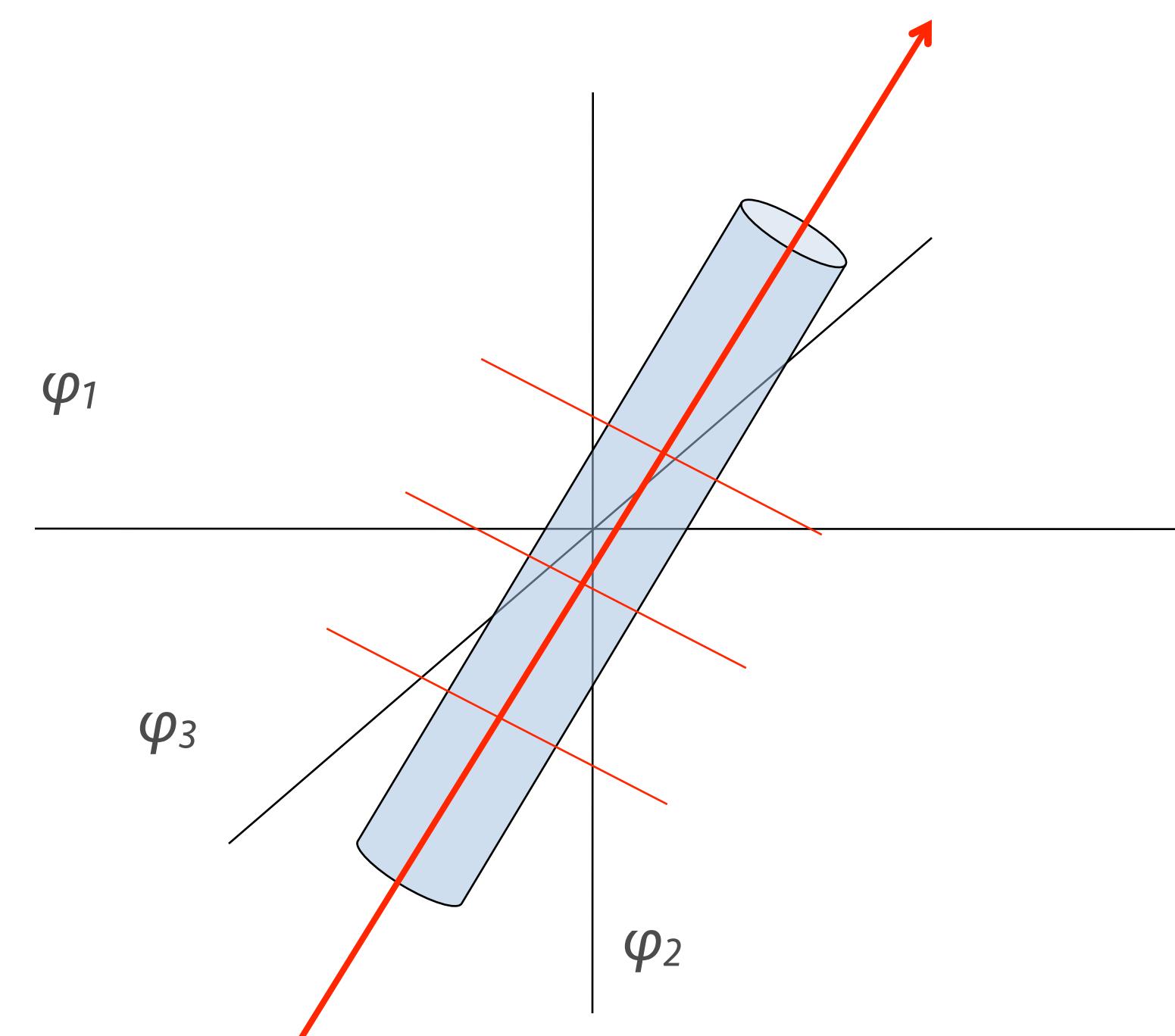


# Now we can do PCA!

- With the right mapping, the new data can be well-behaved for PCA
- The principal components will be in the (potentially higher)  $M$  dimensional space

# Going back to the original space

$$\phi(\mathbf{x}) \in \mathbb{R}^3 \rightarrow \mathbf{x} \in \mathbb{R}^2$$



# From PCA to kernel PCA

- In regular PCA we do:

$$\mathbf{S} = \text{Cov}(\mathbf{x}) \quad \leftarrow \text{Data covariance}$$

$$\mathbf{S} \cdot \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad \leftarrow \text{Eigendecomposition}$$

- In kernel PCA we want to do:

$$\mathbf{C} = \text{Cov}\left(\phi(\mathbf{x})\right) \quad \leftarrow \text{Transformed data covariance}$$

$$\mathbf{C} \cdot \mathbf{v}_i = \lambda_i \mathbf{v}_i \quad \leftarrow \text{Eigendecomposition}$$

# Some problems

- How do we choose the map?
  - We have to operate in this new domain now
- We are in potentially more dimensions
  - More computations == slower!

# Some reshuffling

- From the definition of the eigendecomposition

$$\mathbf{C} \cdot \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

- Each eigenvector  $v_i$  can be described by a linear combination of the input data, so we rewrite the above:

$$\mathbf{v}_i = \sum_n a_{i,n} \phi(\mathbf{x}_n)$$

$$\mathbf{C} \cdot \mathbf{v}_i = \lambda_i \mathbf{v}_i \Rightarrow \frac{1}{N} \sum_n \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \cdot \left( \sum_m a_{im} \phi(\mathbf{x}_m) \right) = \lambda_i \left( \sum_m a_{im} \phi(\mathbf{x}_m) \right)$$

# The kernel trick

- Under certain constraints (more later)

$$\phi(\mathbf{x})^\top \cdot \phi(\mathbf{y}) = K(\mathbf{x}, \mathbf{y})$$

- Where  $K()$  is easier to compute, e.g.:

$$\mathbf{x} = [x_1, x_2]^\top \in \mathbb{R}^2, \quad \phi(\mathbf{x}) = [x_1^2, x_1 x_2, x_2 x_1, x_2^2]^\top \in \mathbb{R}^4$$

$$\phi(\mathbf{x})^\top \cdot \phi(\mathbf{y}) = [x_1^2, x_1 x_2, x_2 x_1, x_2^2] \cdot [y_1^2, y_1 y_2, y_2 y_1, y_2^2]^\top =$$

$$= \left( [x_1, x_2] \cdot [y_1, y_2]^\top \right)^2 = K(\mathbf{x}, \mathbf{y})$$

*4-dimensional map*

*2-dimensional computation!*

# Rewriting the analysis

- The starting high-dim eigendecomposition was:

$$\mathbf{C} \cdot \mathbf{v}_i = \lambda_i \mathbf{v}_i \Rightarrow \frac{1}{N} \sum_n \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^\top \cdot \left( \sum_m a_{im} \phi(\mathbf{x}_m) \right) = \lambda_i \left( \sum_m a_{im} \phi(\mathbf{x}_m) \right)$$

- We can now rewrite as a lower-dim version:

- Multiply on both sides with  $\phi(\mathbf{x}_l)$

$$\begin{aligned} \frac{1}{N} \sum_n K(\mathbf{x}_l, \mathbf{x}_n) \sum_m a_{im} K(\mathbf{x}_l, \mathbf{x}_m) &= \lambda_i \sum_m a_{im} K(\mathbf{x}_l, \mathbf{x}_m) \\ \Rightarrow \mathbf{K}^2 \cdot \mathbf{a}_i &= \lambda_i N \mathbf{K} \cdot \mathbf{a}_i \end{aligned}$$

# That's solved easily

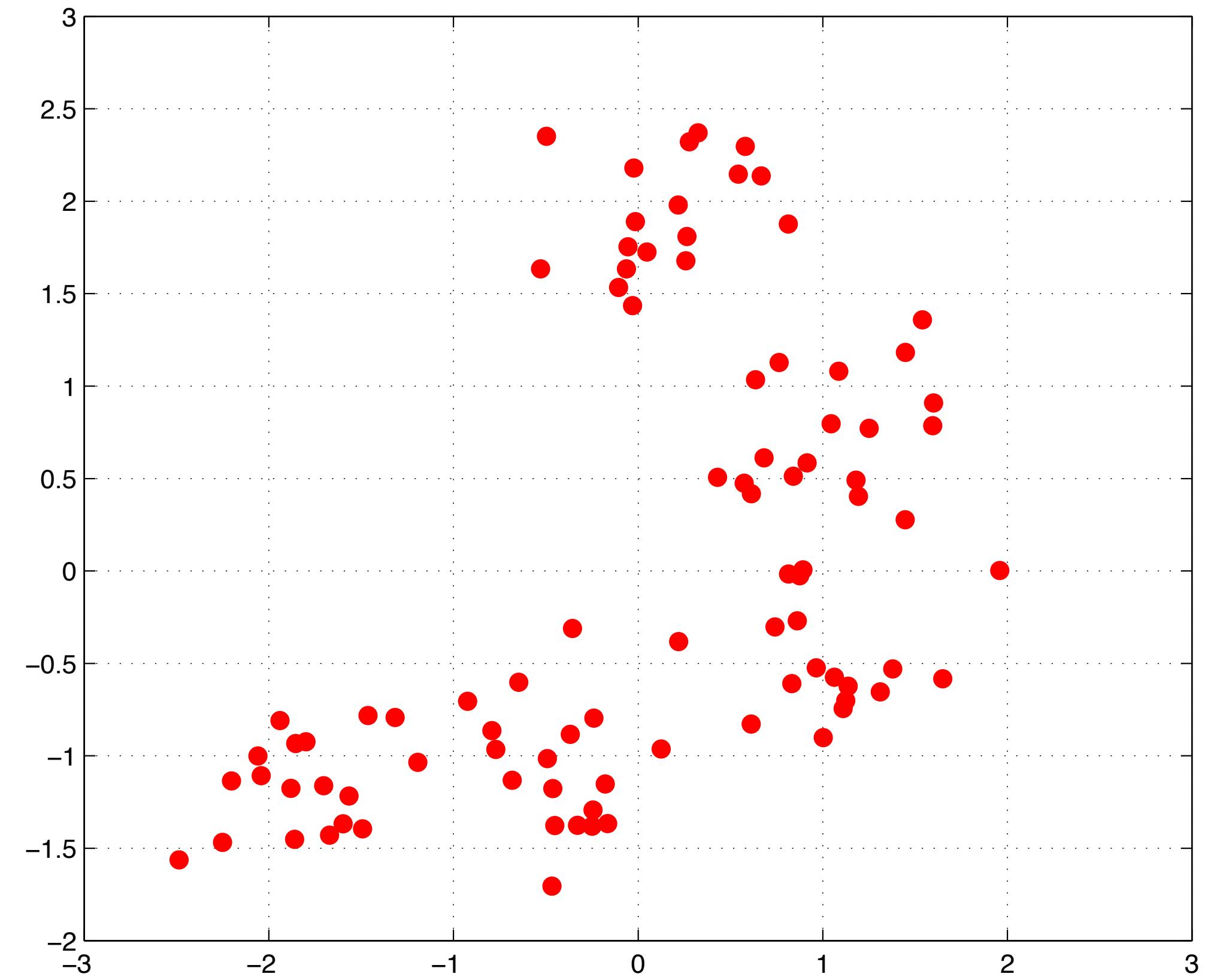
- It is an eigendecomposition once again!

$$\mathbf{K} \cdot \mathbf{a}_i = \lambda_i N \mathbf{a}_i$$

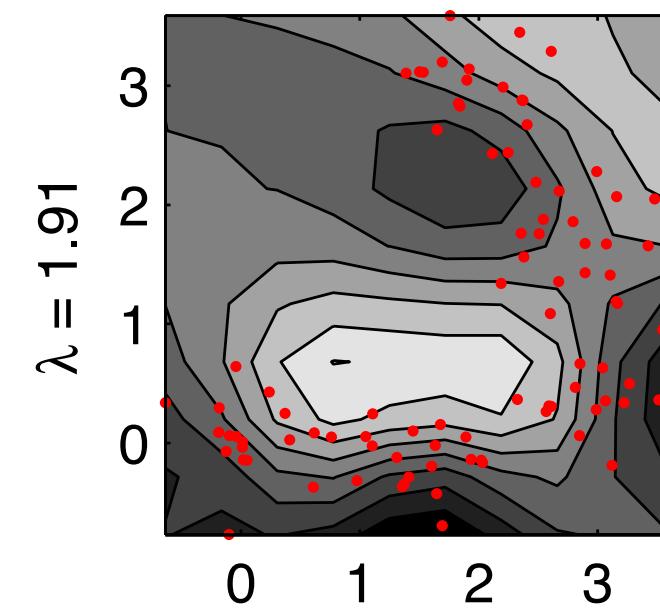
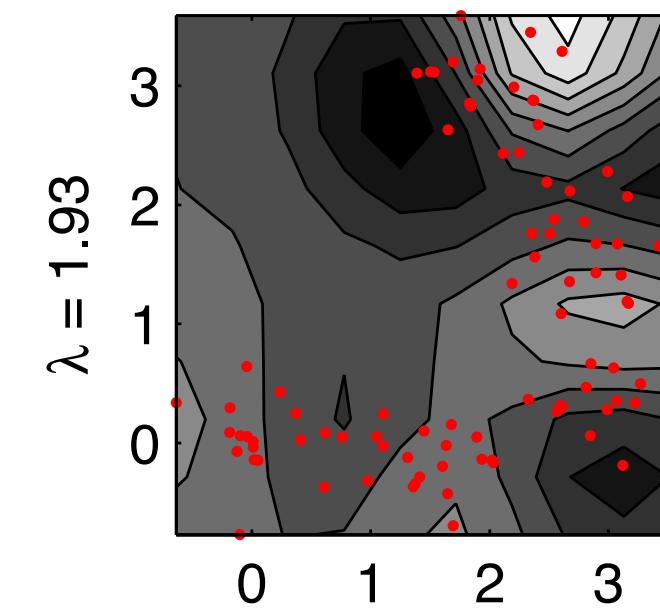
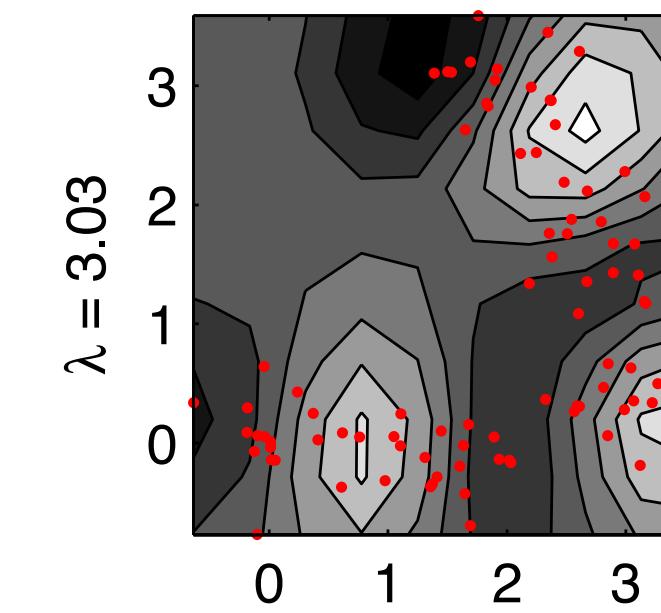
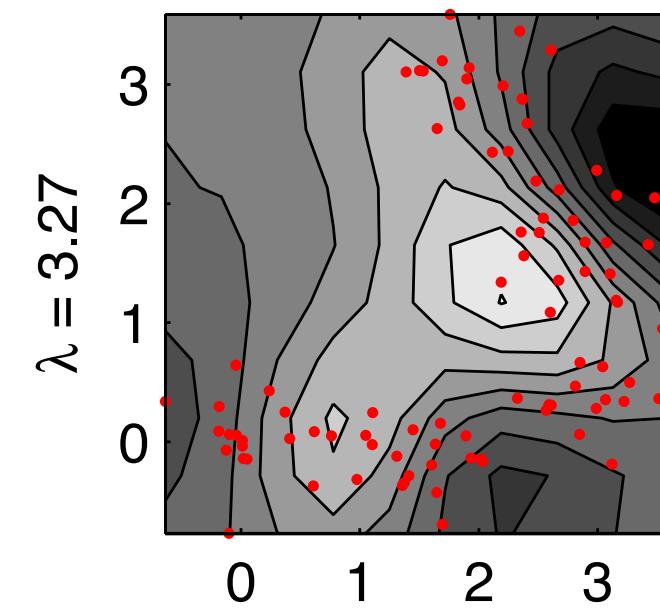
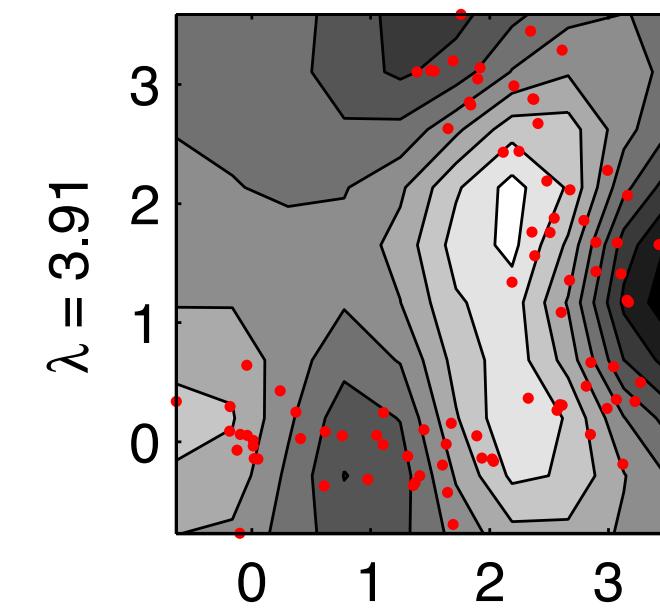
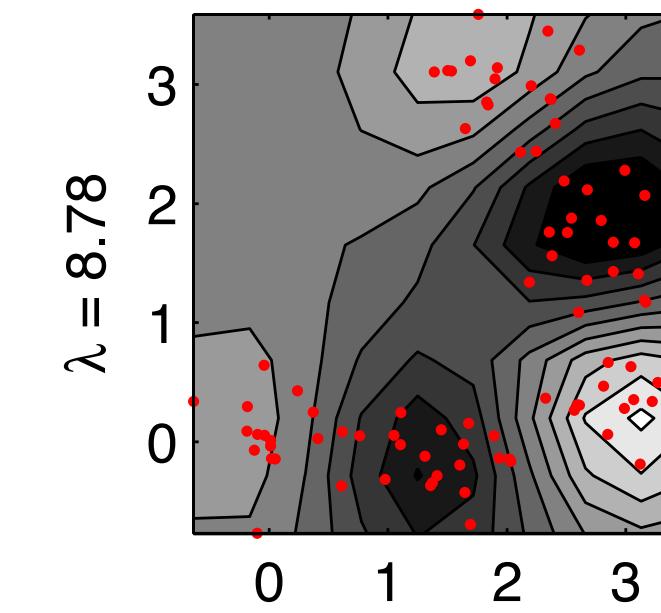
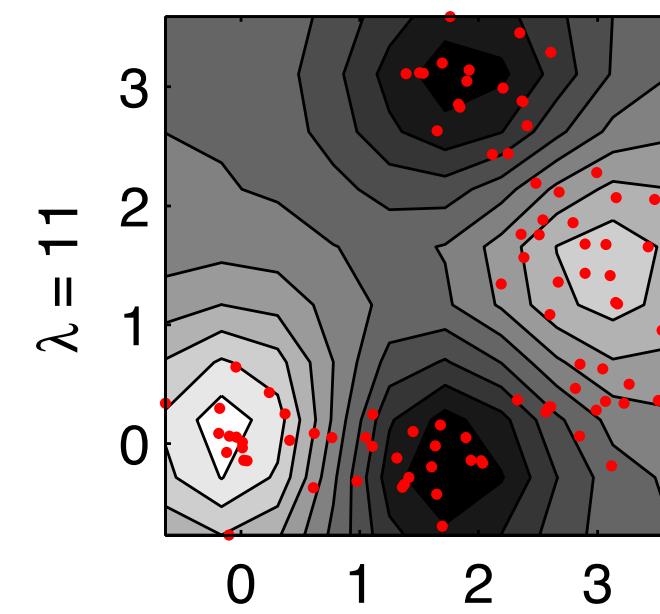
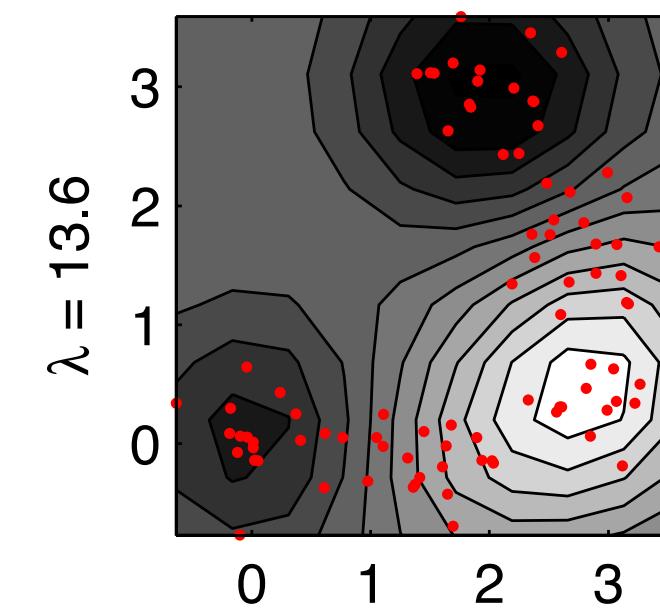
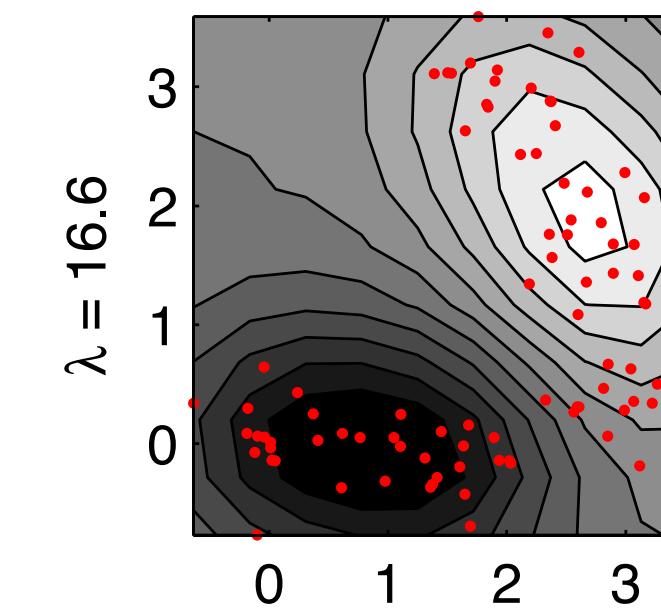
- The non-linear principal components are  $a_i$
- The matrix  $\mathbf{K}$  is now  $N \times N$
- In practice we need to ensure that  $\mathbf{K}$  is also positive definite (i.e.  $\varphi(\mathbf{x})$  is zero mean)
  - Not hard to do, but we skip the details for now

# An example

- Data is not good for PCA
  - too curvy, not Gaussian
- Use this kernel:
$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2}$$
- Eigenvectors will be in terms of some high-D space



# The non-linear eigenvectors



# Some kPCA notes

- The eigenanalysis is now bigger
  - $\mathbf{K}$  is  $N \times N$ 
    - Instead of  $\mathbf{X} \cdot \mathbf{X}^T$  we analyze  $\varphi(\mathbf{X}^T \cdot \mathbf{X})$
    - So computations will be slower
  - How many components can we extract?
    - Not more than the original dimensions
      - There will be zero eigenvalues past that
  - How do we choose the right kernel? Who knows ...

# A different take on PCA

- What happens when we only have quantified relationships between data?
  - E.g. distance or similarity metrics
- Can we do dimensionality reduction?
- How do we find the geometry of our data?
  - Very useful in psychology and social sciences

# Reconstructing from distances

- Suppose we only have cross-data distances
  - Euclidean ones for now

- E.g.:

*Distance matrix*

	$x_1$	$x_2$	$x_3$
$x_1$	0	1	1
$x_2$	1	0	1
$x_3$	1	1	0

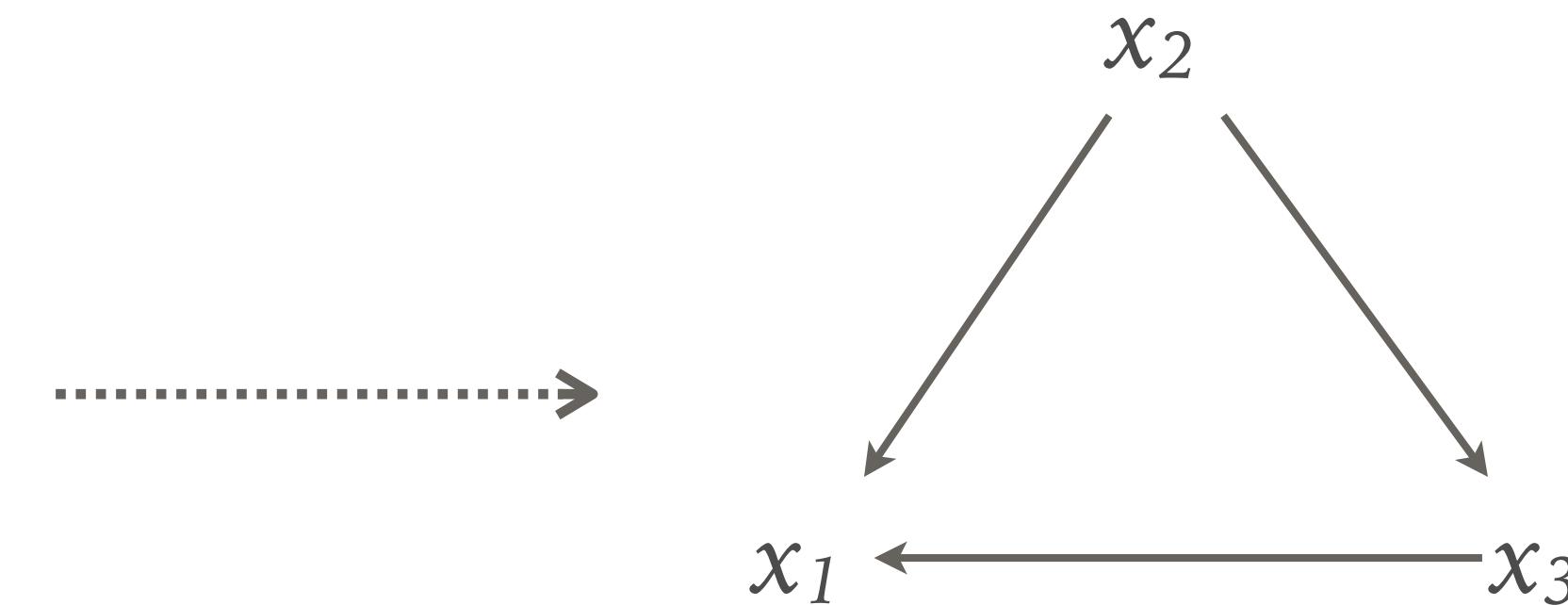
# Reconstructing from distances

- Suppose we only have cross-data distances
  - Euclidean ones for now

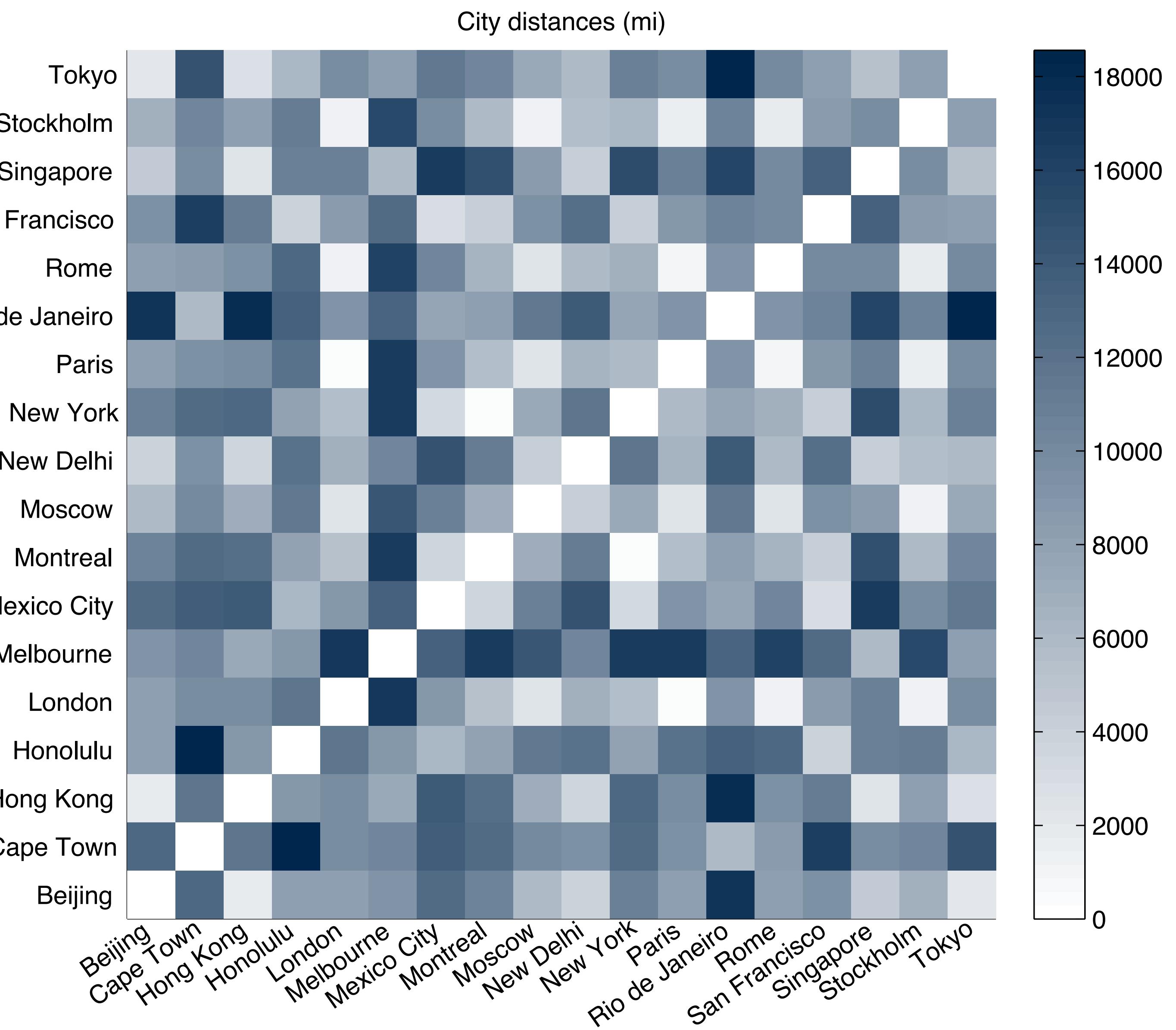
- E.g.:

*Distance matrix*

	$x_1$	$x_2$	$x_3$
$x_1$	0	1	1
$x_2$	1	0	1
$x_3$	1	1	0



# What about now?



# Getting from distances to points

- Given only the cross-point distances:

$$d_{i,j} = \left\| \mathbf{x}_i - \mathbf{x}_j \right\|^2$$

- We want to estimate  $\mathbf{x}$ , such that:

$$\left\| \hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j \right\|^2 \approx d_{i,j}$$

- How do we go about it?

# Writing this as a matrix operation

- We can express the squared distances as a product:

$$d_{i,j} = \left\| \mathbf{x}_i - \mathbf{x}_j \right\|^2 = \mathbf{x}_i^\top \cdot \mathbf{x}_i + \mathbf{x}_j^\top \cdot \mathbf{x}_j - 2\mathbf{x}_i^\top \cdot \mathbf{x}_j \Rightarrow$$
$$\mathbf{D} = \text{diag}\left(\mathbf{X}^\top \cdot \mathbf{X}\right) \cdot \mathbf{1}_N^\top + \mathbf{1}_N \cdot \text{diag}\left(\mathbf{X}^\top \cdot \mathbf{X}\right)^\top - 2\mathbf{X}^\top \cdot \mathbf{X}$$



- If the first two terms were zero, we could solve:

$$\mathbf{D} = \mathbf{X}^\top \cdot \mathbf{X}$$

- But they are not, so how do we remove the unwanted terms?
  - Let's average them out

# Removing the unwanted terms

- To remove the row/column mean of a matrix do:

$$\mathbf{P} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \cdot \mathbf{1}_N^\top = \mathbf{I}_N - \frac{\mathbf{1}_{N \times N}}{N}$$

$$\underbrace{(\mathbf{A} \cdot \mathbf{P}) \cdot \frac{1}{N} \mathbf{1}_N}_{\text{Remove mean column}} = \mathbf{0}_N \quad \frac{1}{N} \mathbf{1}_N^\top \cdot \underbrace{(\mathbf{P} \cdot \mathbf{A})}_{\text{Remove mean row}} = \mathbf{0}_N^\top$$

- Applying on  $\mathbf{D}$ :

$$\begin{aligned} \mathbf{P} \cdot \mathbf{D} \cdot \mathbf{P} &= \mathbf{P} \cdot \text{diag}(\mathbf{X}^\top \cdot \mathbf{X}) \cdot \mathbf{1}_N^\top \cdot \mathbf{P} + \mathbf{P} \cdot \mathbf{1}_N \cdot \text{diag}(\mathbf{X}^\top \cdot \mathbf{X}) \cdot \mathbf{P} - 2\mathbf{P} \cdot \mathbf{X}^\top \cdot \mathbf{X} \cdot \mathbf{P} = \\ &= -2\mathbf{P} \cdot \mathbf{X}^\top \cdot \mathbf{X} \cdot \mathbf{P} = -2(\mathbf{X} \cdot \mathbf{P})^\top \cdot (\mathbf{X} \cdot \mathbf{P}) = -2(\mathbf{X} - \bar{\mathbf{x}})^\top \cdot (\mathbf{X} - \bar{\mathbf{x}}) \end{aligned}$$

- Bonus, a zero mean assumption on  $\mathbf{x}$  (which we like!)

# Overall process

- Get distances and “double center”:

$$S = -\frac{1}{2} \left( D - \frac{1}{N} D \cdot \mathbf{1}_N \cdot \mathbf{1}_N^\top - \frac{1}{N} \mathbf{1}_N \cdot \mathbf{1}_N^\top \cdot D + \frac{1}{N^2} \mathbf{1}_N \cdot \mathbf{1}_N^\top \cdot D \cdot \mathbf{1}_N \cdot \mathbf{1}_N^\top \right)$$

- Estimate  $\mathbf{X}$  from  $S$

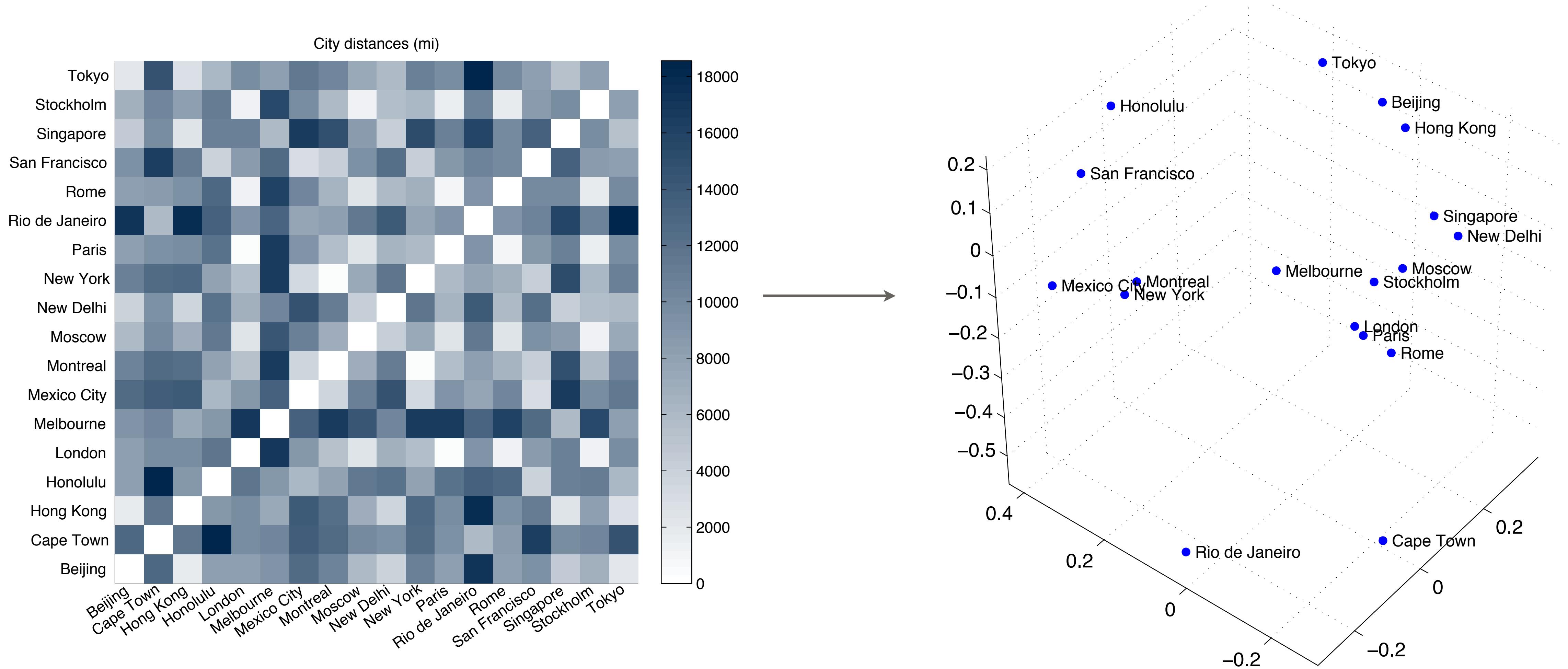
$$S = \mathbf{X}^\top \cdot \mathbf{X} \Rightarrow S = \mathbf{U} \cdot \Lambda \cdot \mathbf{U}^\top \Rightarrow \hat{\mathbf{X}} = \Lambda^{1/2} \cdot \mathbf{U}^\top$$

Eigendecomposition

- $\mathbf{X}$  is an  $N$ -dimensional representation of the data points
  - Keep as many dimensions you want
    - They are ordered by importance from the eigendecomposition!

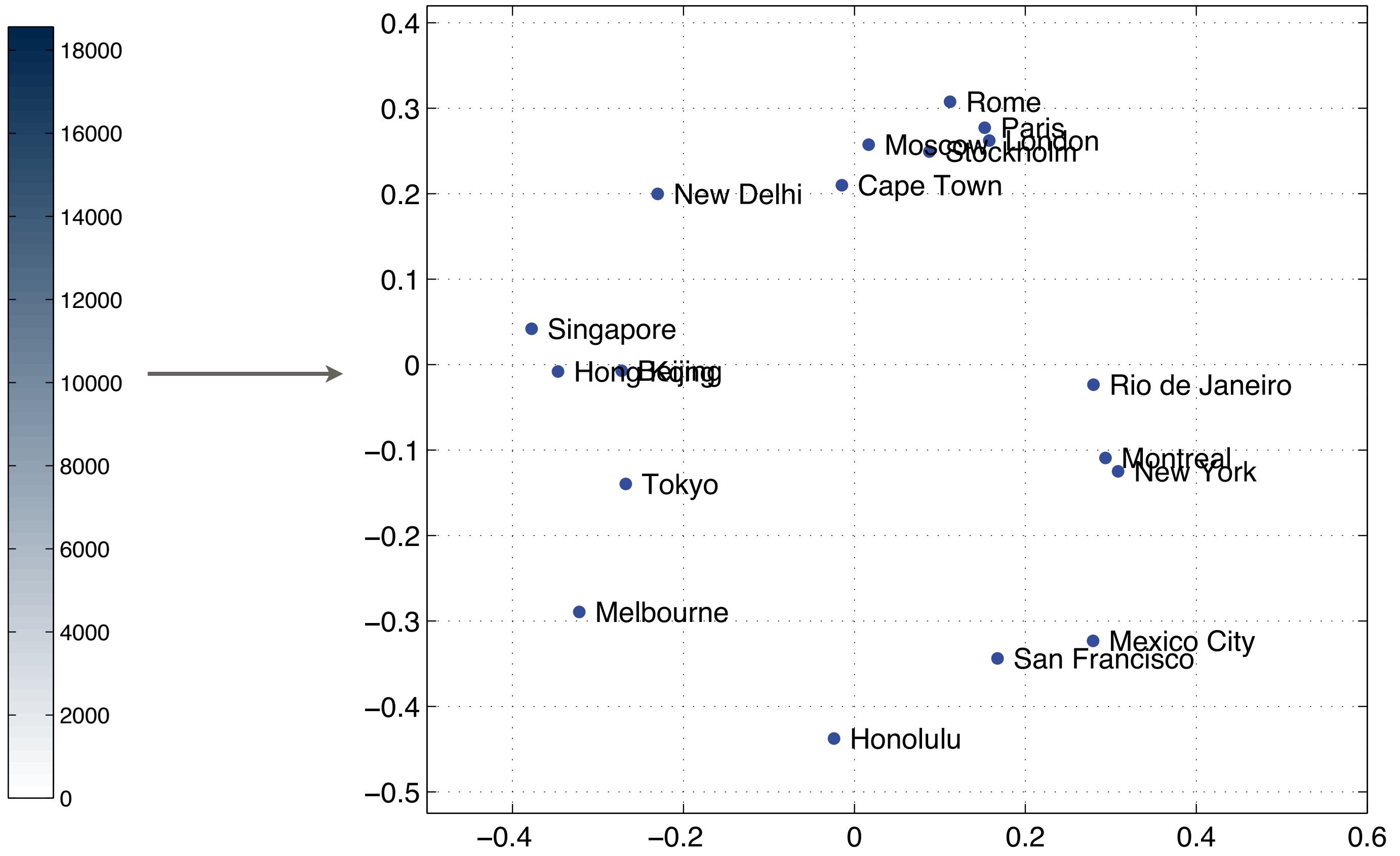
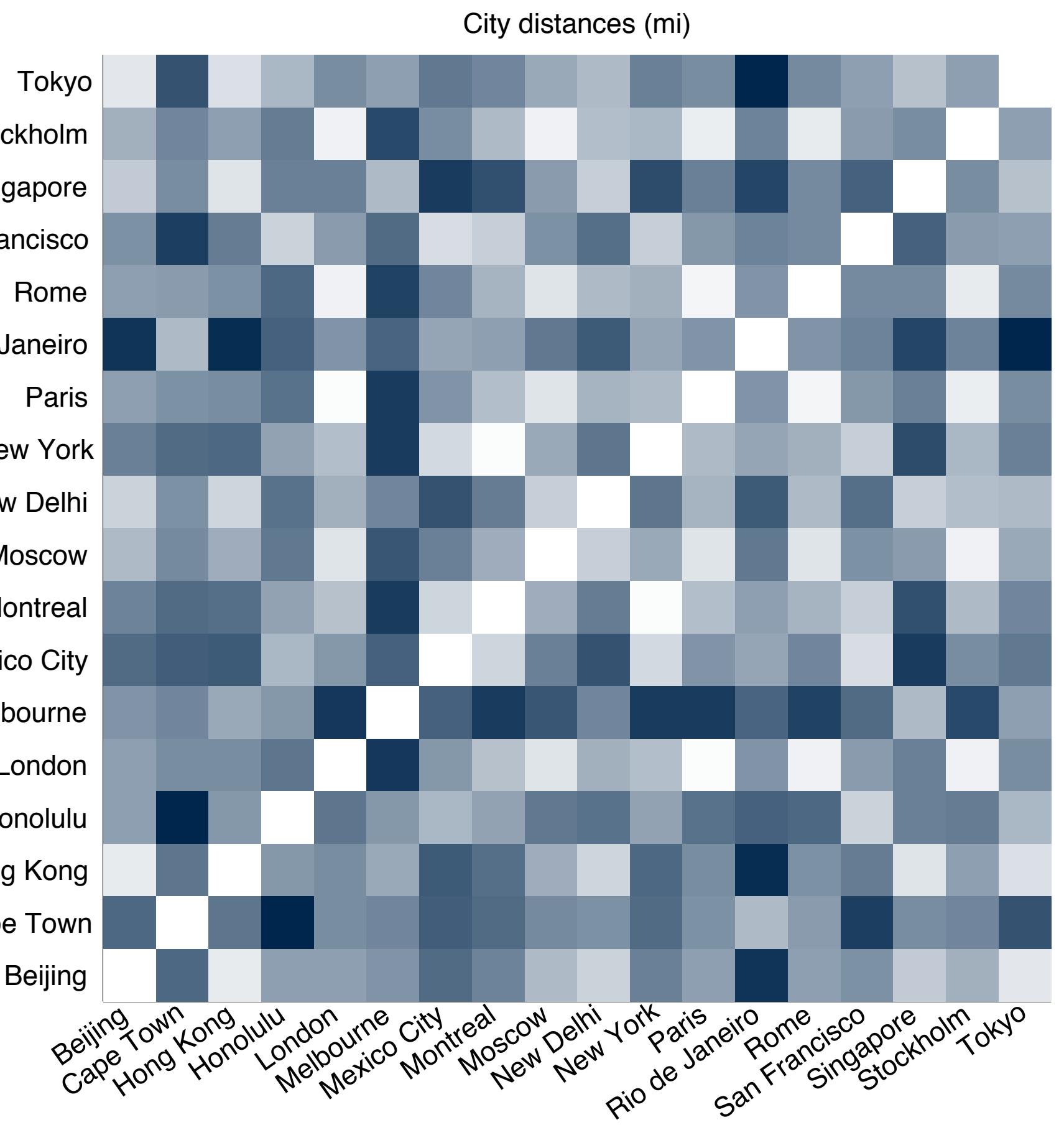
# Back to the problem at hand

- City distance matrix results in an accurate map!



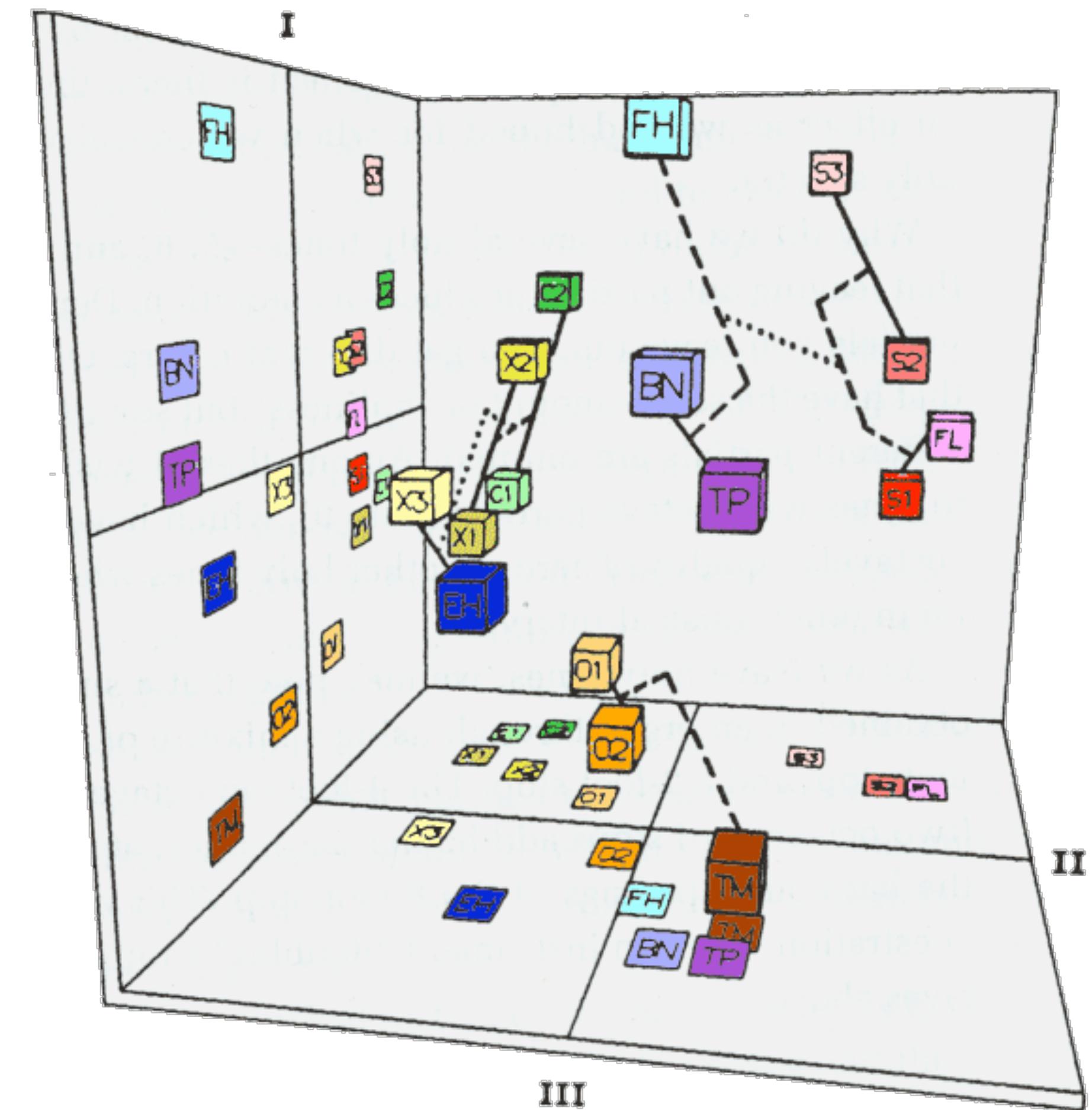
# With fewer dimensions

- Also pretty accurate!



# Another use

- Interpreting percepts
  - Get distances from human subjects
  - Create map of percept
- Gray's timbre space



# Variations

- Preserving dot products
- Using ranking as opposed to distances
- Using non-linear distance metrics
- Using only local data

# What is MDS good for

- Finds the implied geometry of our data
  - Not obvious for large dimensions!
- It is a sort of like PCA
  - Like PCA: But decomposes  $\mathbf{X}^\top \cdot \mathbf{X}$  not  $\mathbf{X} \cdot \mathbf{X}^\top$
  - Like kPCA: The kernel is the inner product
  - Is linear

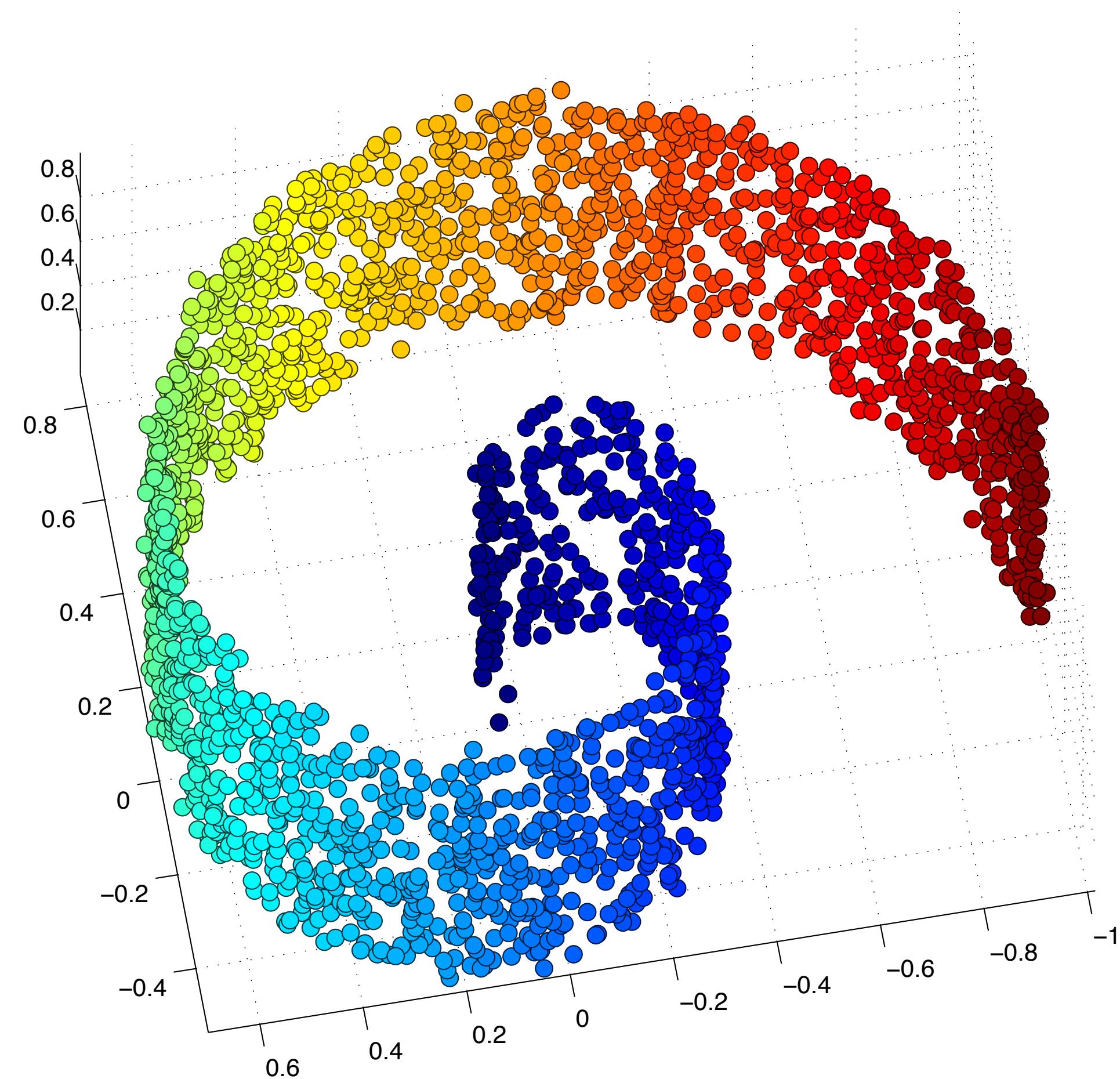
# Data with low intrinsic dimension

- Actual dimension can be misleading



# Data with low intrinsic dimension

- Actual dimension can be misleading



# Local vs. Global



# Local vs. Global



# Local vs. Global



# Preserving topology

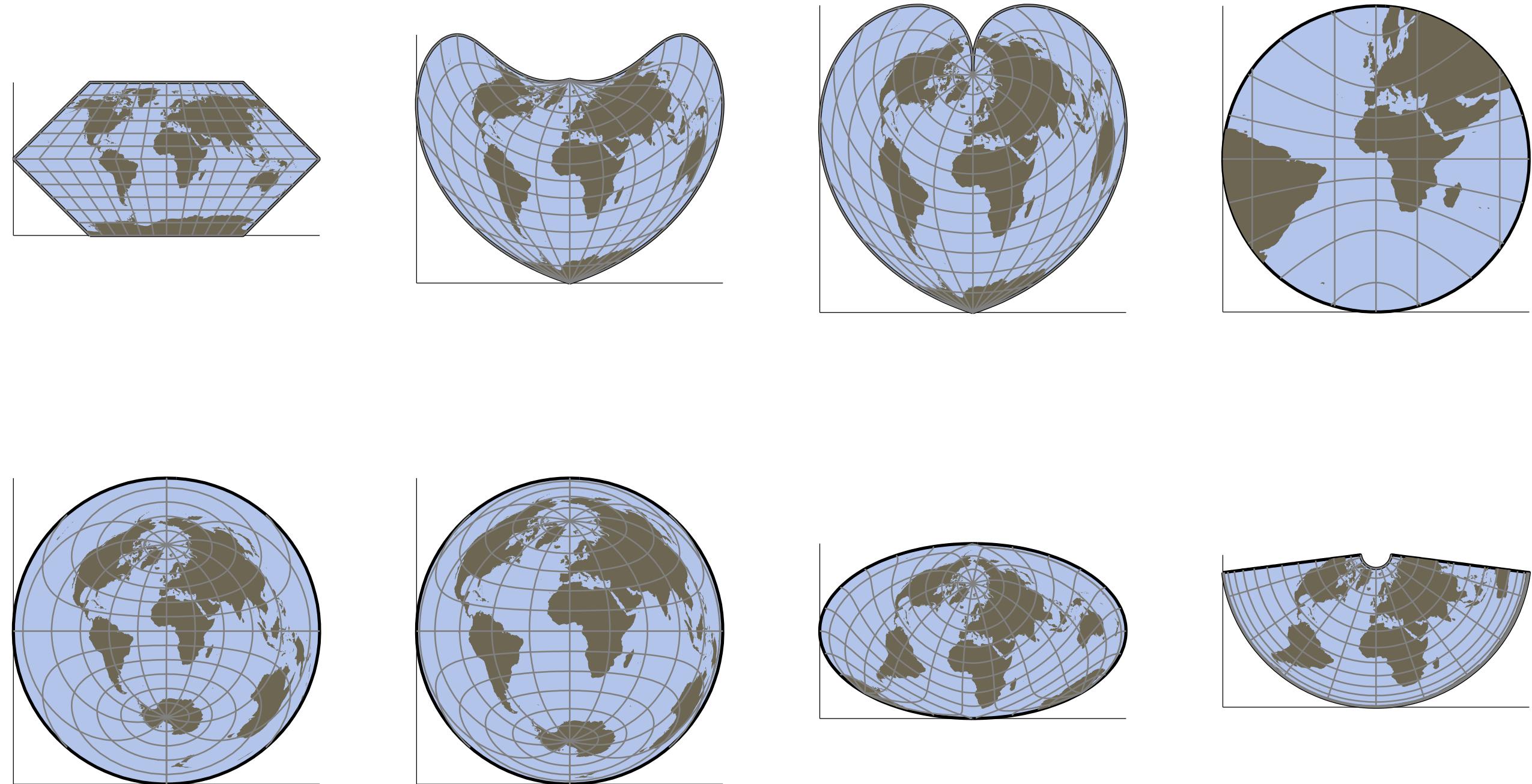
- Looking at the neighboring structure
  - Is the earth flat or round?
- Reduce dimensionality while keeping the neighborhood structure unchanged

# Example: Map projections

- Many ways to go from 3D to 2D
  - Depends on what you care about
    - How do we do this for data?

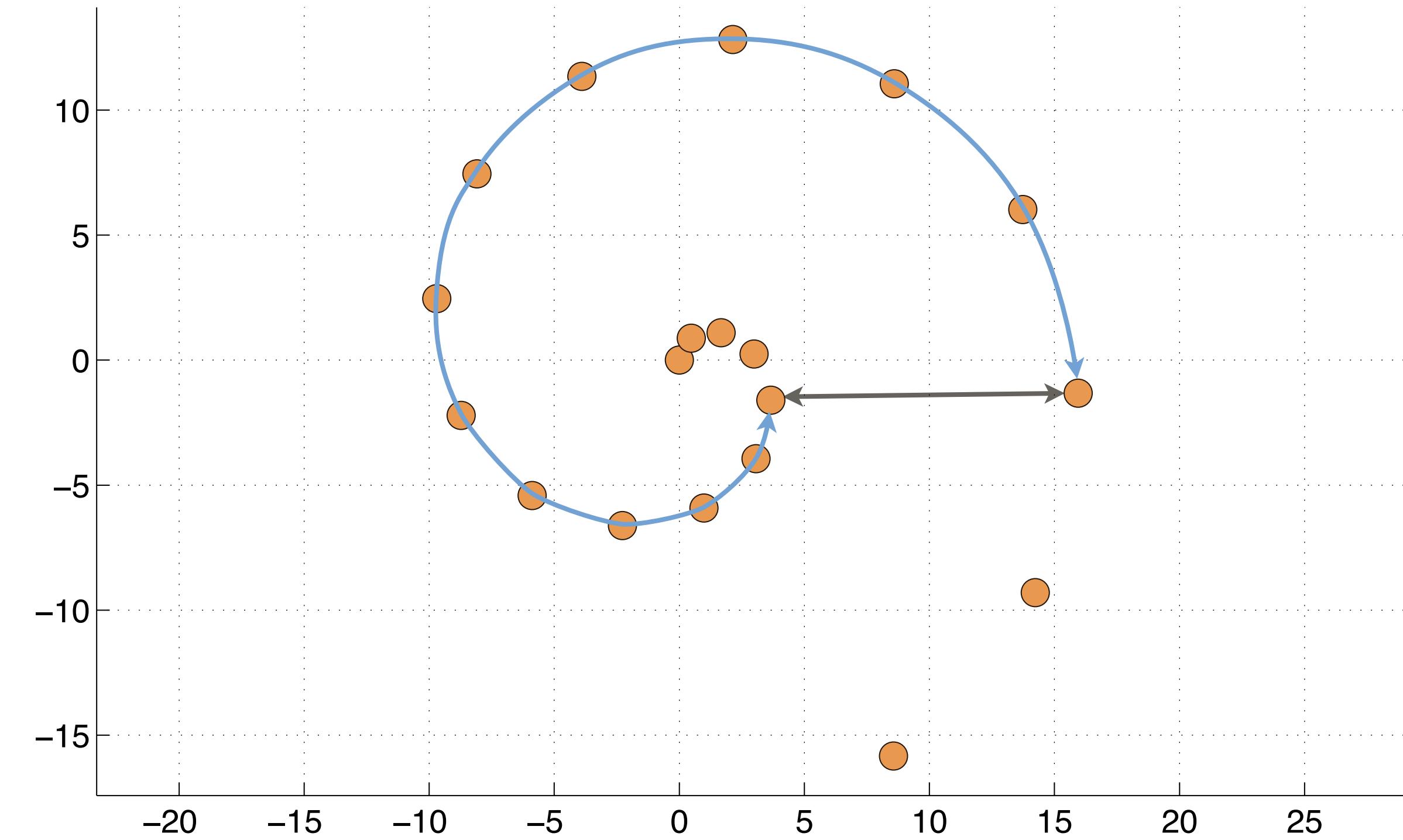


*Dimensionality  
reduction*



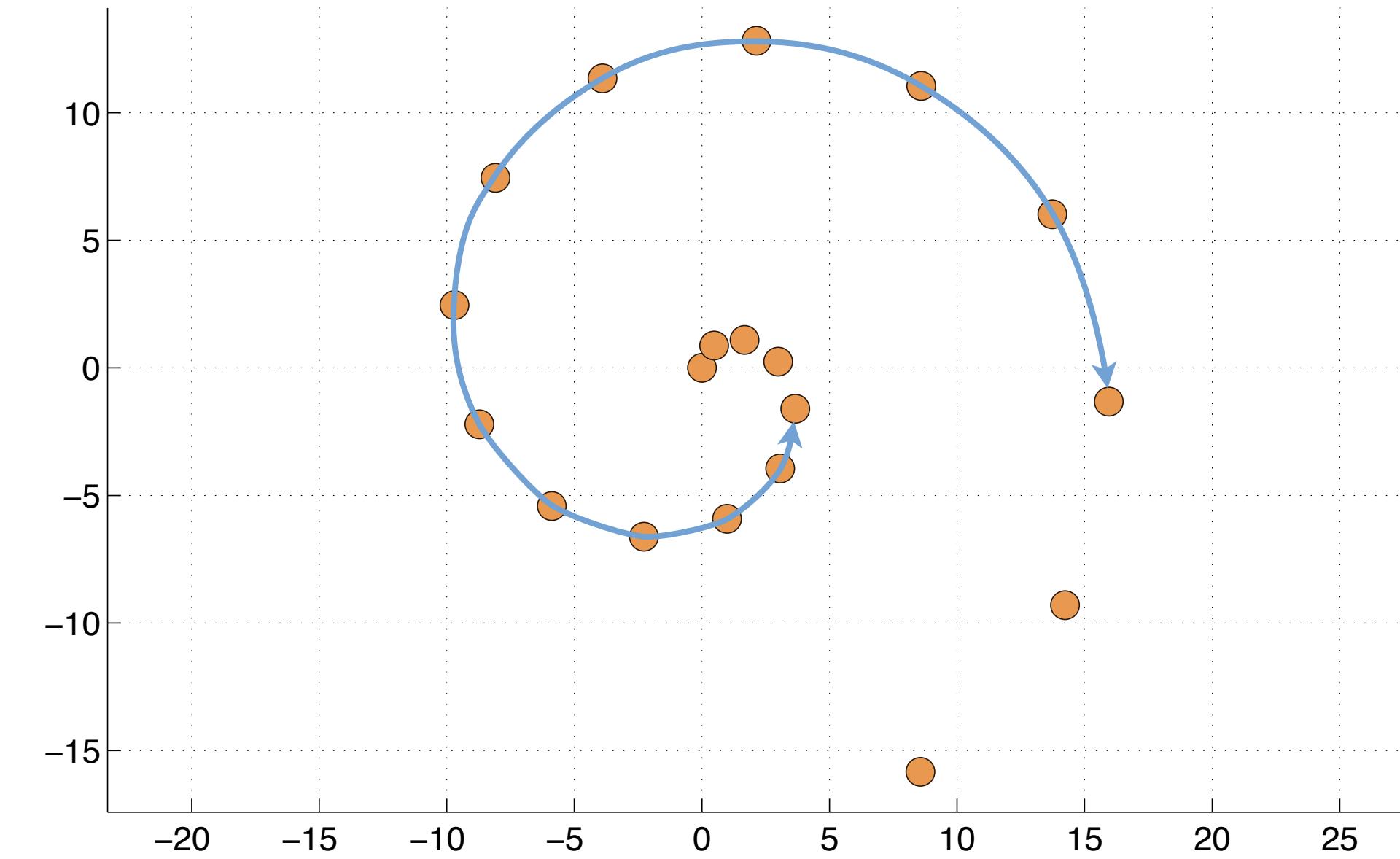
# Euclidean vs. Geodesic distance

- On a manifold distance is subjective
  - Euclidean distance ignores manifold structure
  - Geodesic distance is more appropriate

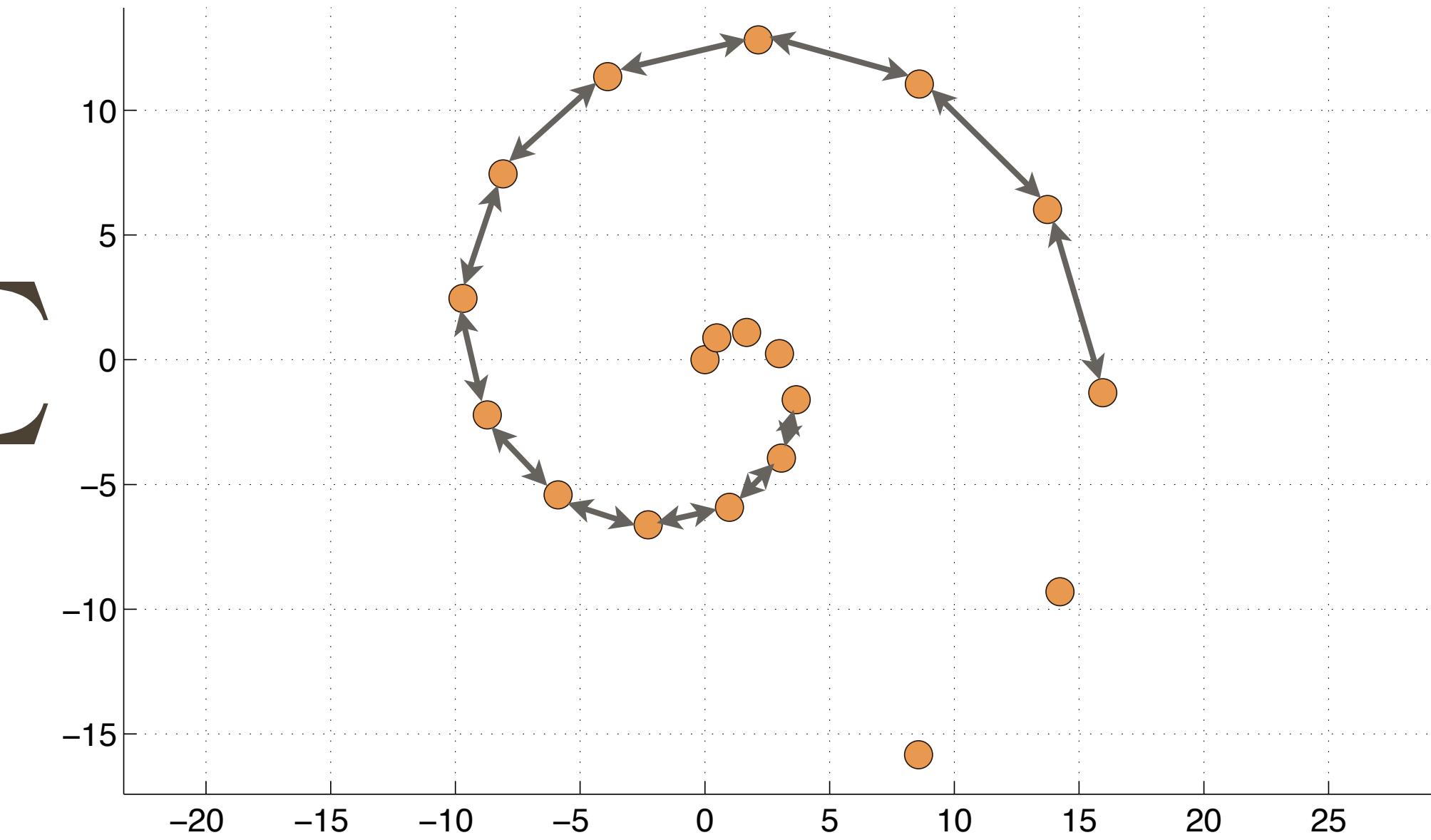


# Getting the distance using just data

- Since we don't know the manifold we can use the available data to get the geodesic distance
  - Shortest path between two points that passes through neighboring data



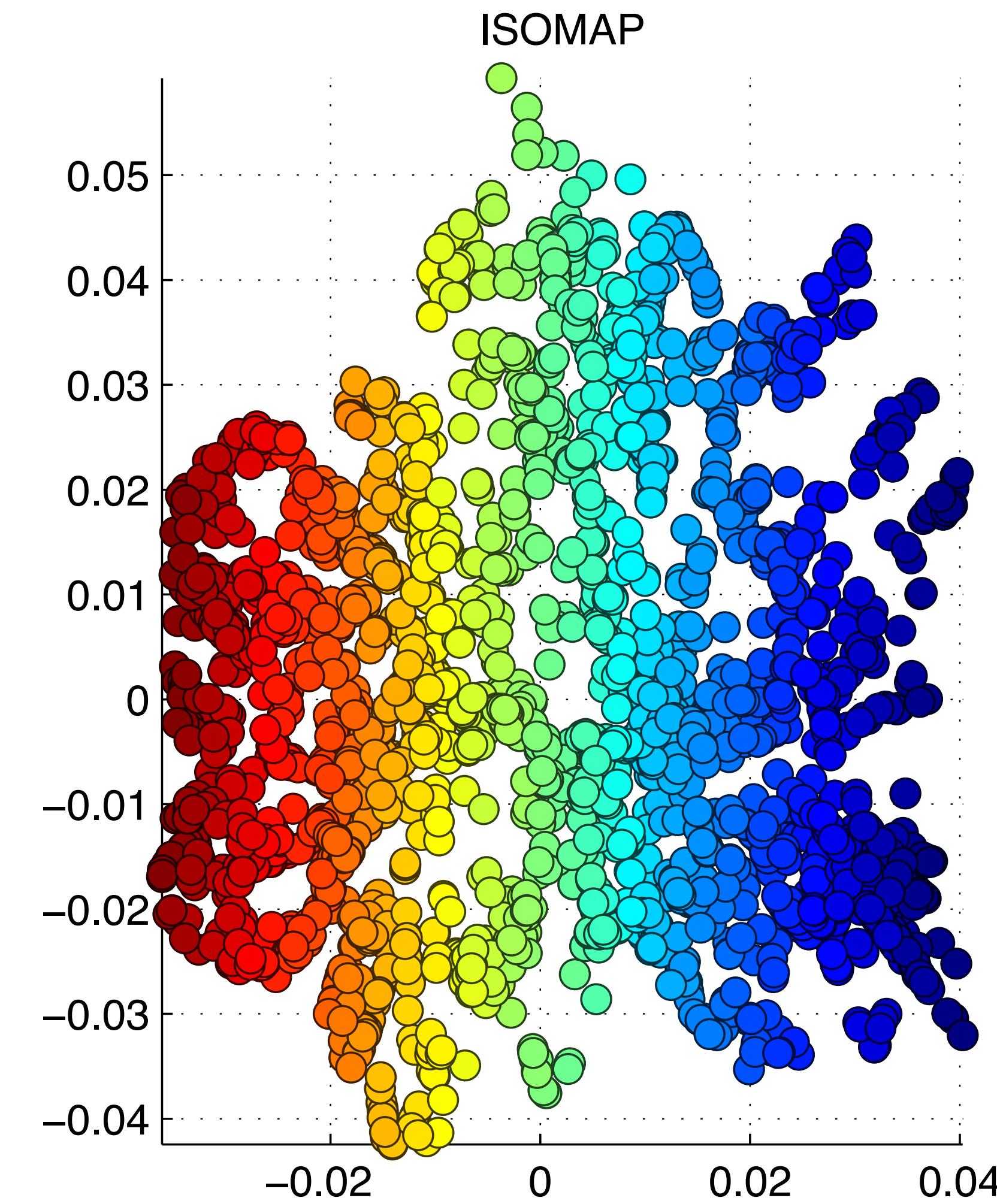
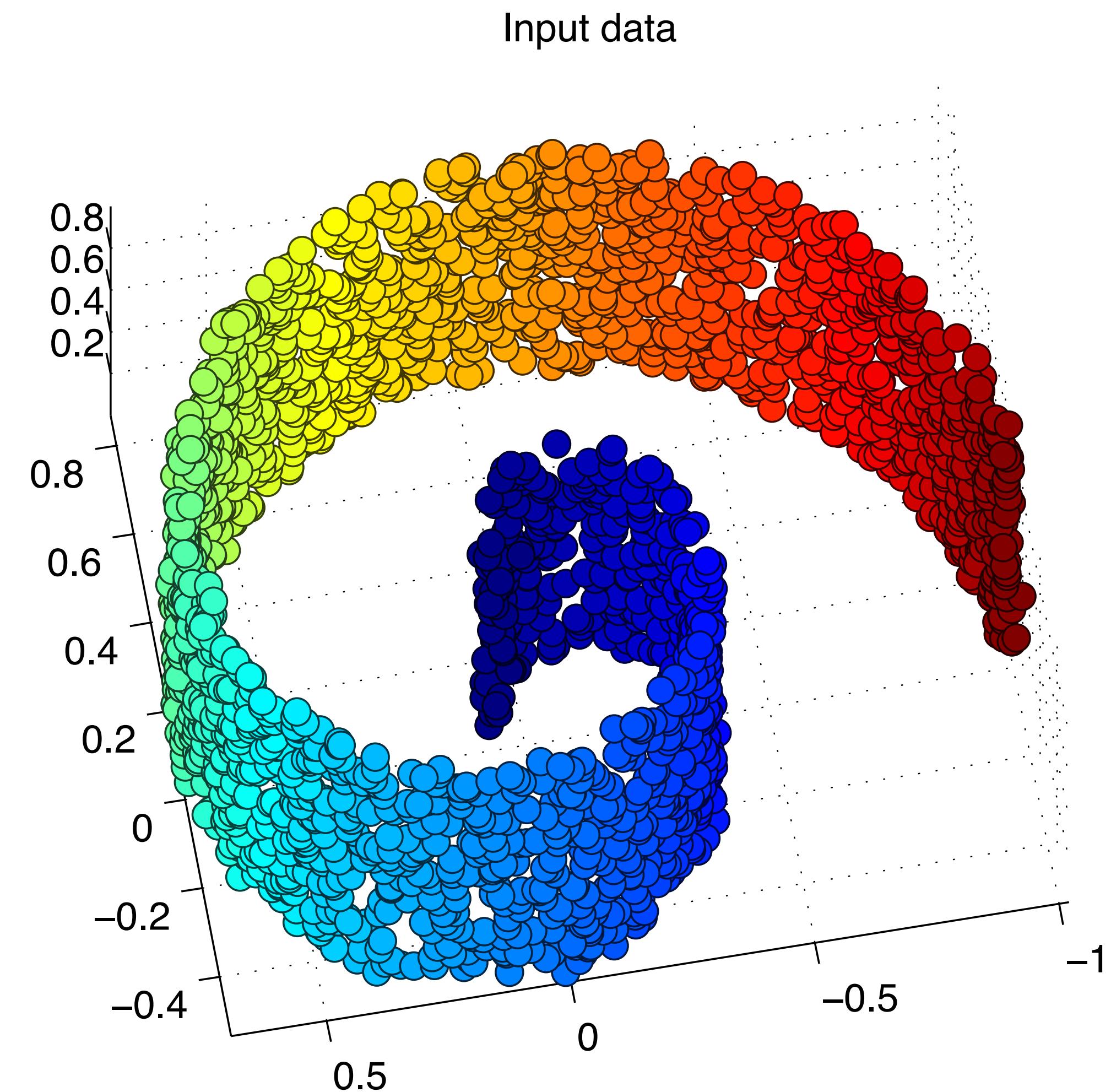
$$\approx \sum$$



# ISOMAP

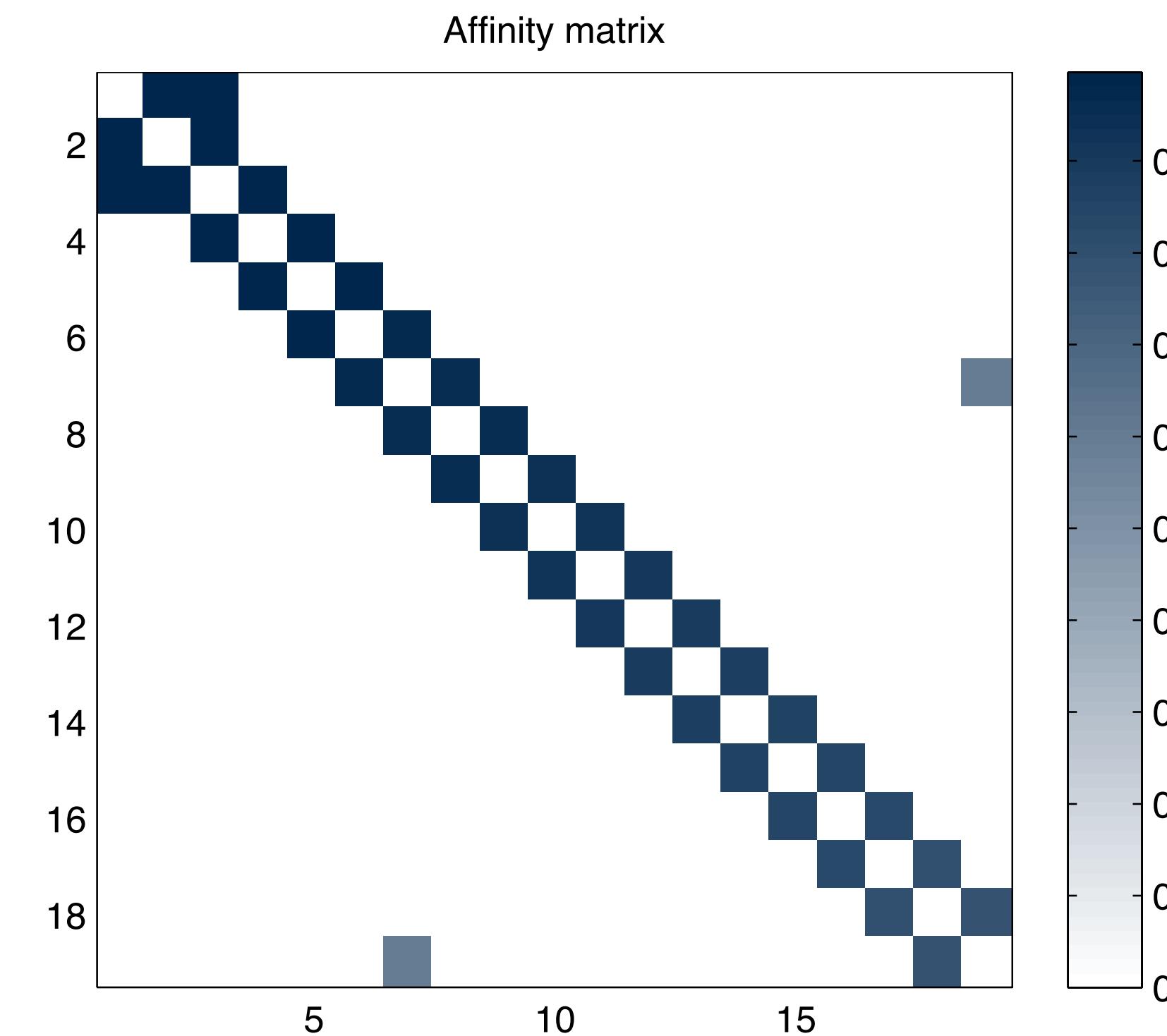
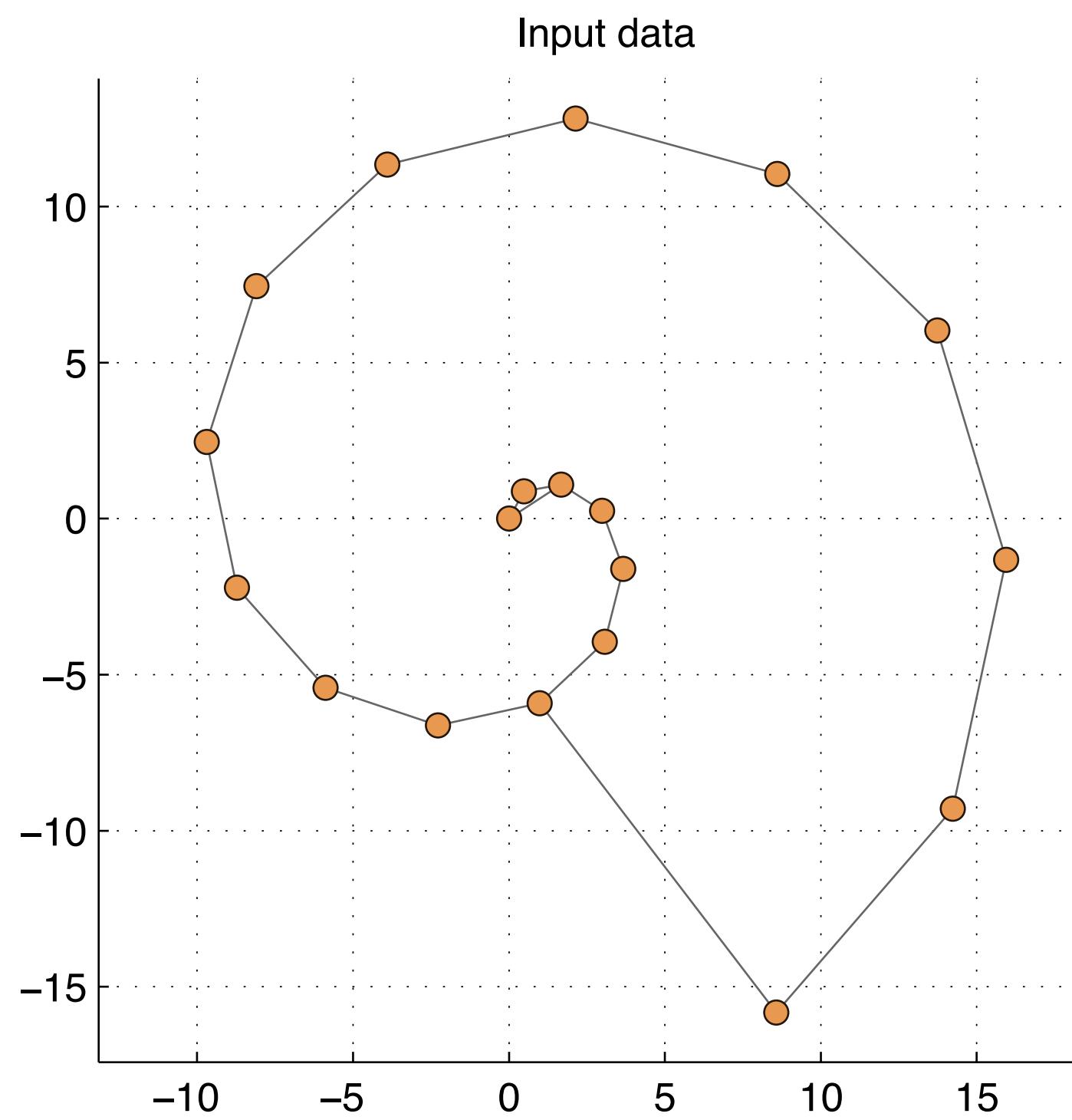
- Perform MDS using geodesic distances
  - i.e. change only  $D$ , the matrix with the distances
- Results in an embedding that's aware of the low-dimensional structure in our data
  - Results in Euclidean distances that are approximately equal to the geodesic distances

# On the Swiss roll



# A neighborhood graph approach

- Make a note of point distances
  - $N$ -nearby points matter, zero out the rest (and the diagonal)



$$w_{i,j} = e^{-\frac{\|x_i - x_j\|^2}{\sigma}}$$

# A goal

- Get an embedding that minimizes:

$$E_{LE} = \sum_{i,j} \left\| \mathbf{y}_i - \mathbf{y}_j \right\|^2 w_{i,j}$$

- We can rewrite as:

$$E_{LE} = 2\mathbf{Y}^\top \cdot \mathbf{L} \cdot \mathbf{Y}$$

$$\mathbf{L} = \mathbf{W} - \text{diag}(\mathbf{1}^\top \cdot \mathbf{W}) = \mathbf{W} - \mathbf{D}$$

*Laplacian of the distance graph* → ↑ *Sum of rows*

	Large $ \mathbf{y}_i - \mathbf{y}_j $	Small $ \mathbf{y}_i - \mathbf{y}_j $
Large $w_{ij}$	<b>Bad</b>	Good
Small $w_{ij}$	Good	Don't care

# And we solve

- Impose a constraint against arbitrary scaling:

$$\text{minimize } \mathbf{Y}^\top \cdot \mathbf{L} \cdot \mathbf{Y}$$

$$\text{subject to } \mathbf{Y}^\top \cdot \mathbf{D} \cdot \mathbf{Y} = \mathbf{I}$$

To avoid  $y_i = 0$   
and arbitrary scaling

- one more rewrite:

$$\text{minimize } \mathbf{Z}^\top \cdot \tilde{\mathbf{L}} \cdot \mathbf{Z}$$

$$\text{subject to } \mathbf{Z}^\top \cdot \mathbf{Z} = \mathbf{I}$$

Normalized graph Laplacian

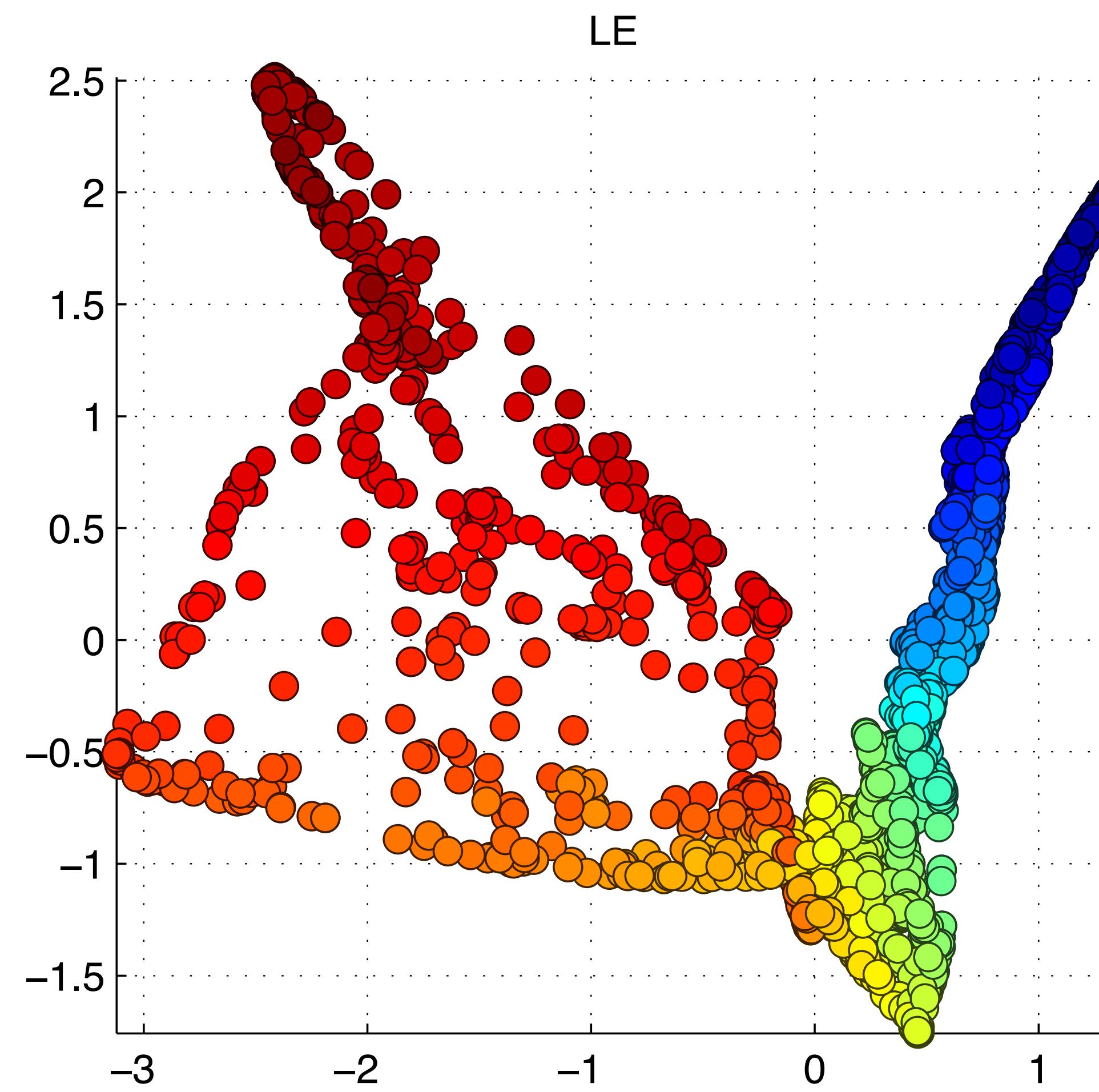
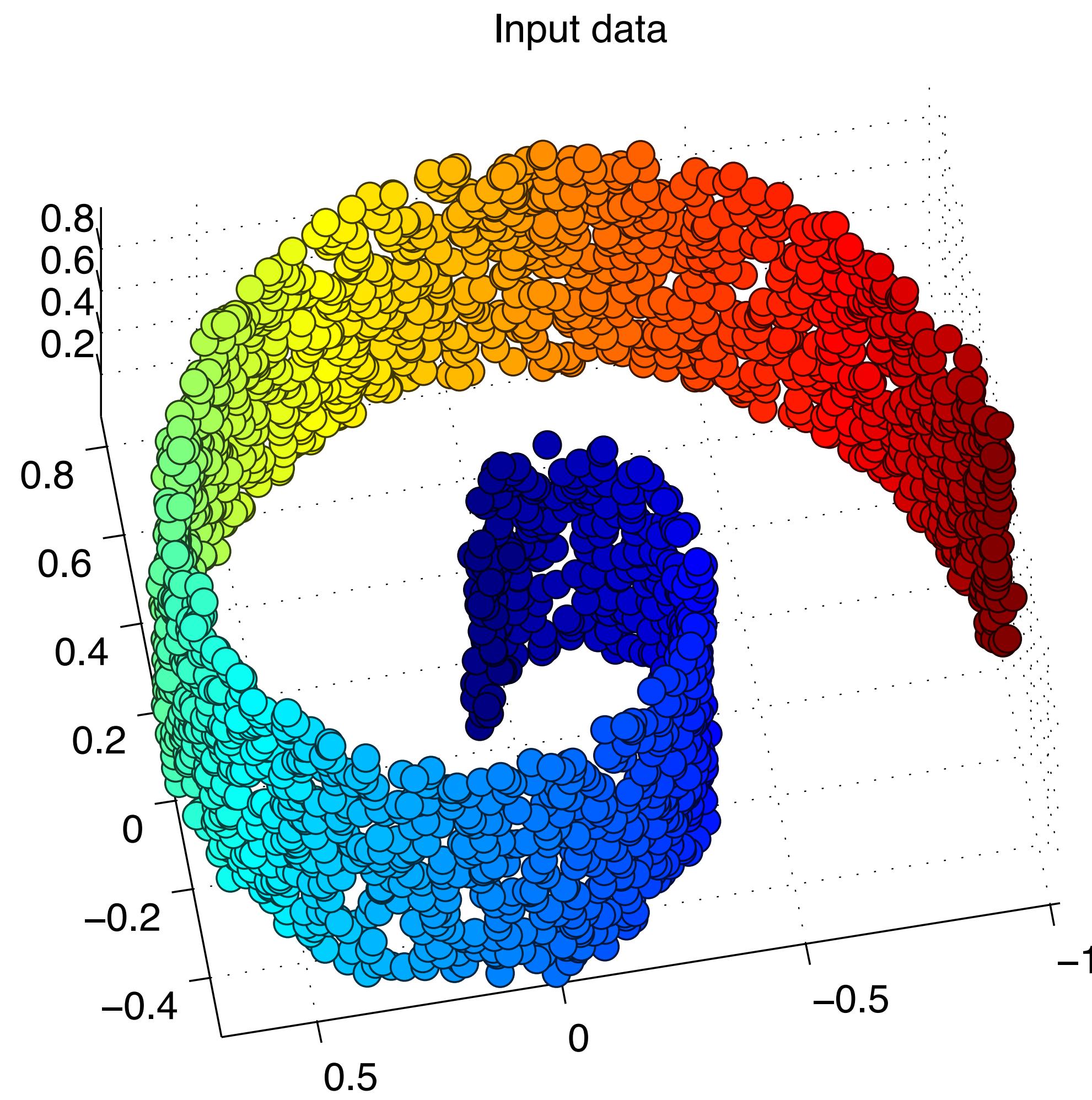
$$\mathbf{Z} = \mathbf{D}^{1/2} \cdot \mathbf{Y}$$

$$\tilde{\mathbf{L}} = \mathbf{D}^{-1/2} \cdot \mathbf{L} \cdot \mathbf{D}^{-1/2}$$

# And predictably ...

- $Z$  are the eigenvectors of the normalized Laplacian
  - From which we get  $Y$ 
$$Y = Z^\top \cdot D^{1/2}$$
- The eigenvalues will help sort  $Y$ 
  - Small eigenvalues imply closer spaced  $y$ 's
    - The smallest eigenvalue will have all  $y$ 's on one point (useless)
- So we pick the  $2^{nd}$  to  $N^{th} + 1$  smallest eigenvectors to produce an embedding for  $N$  dimensions

# Swiss roll example



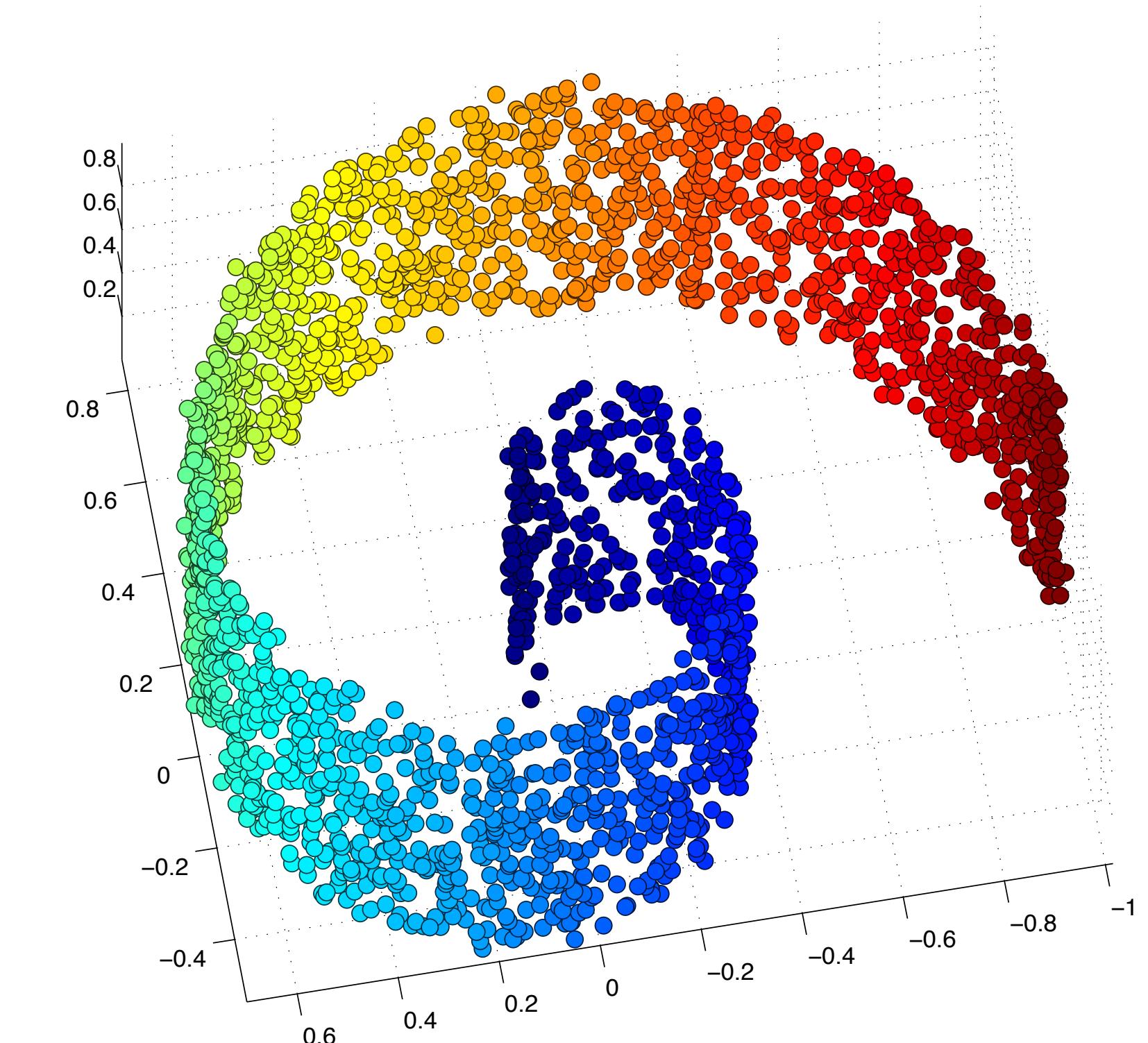
# Locally Linear Embedding

- Observe local neighborhood of all points
  - Assume each neighborhood is linear
  - Explain each point using its neighbors

$$\mathbf{x}_i \approx \sum_{j \in N(i)} w_{i,j} \mathbf{x}_j$$

- And there's an optimal  $\mathbf{W}$

$$C_H(\mathbf{W}) = \sum_i \left\| \mathbf{x}_i - \sum_{j \in N(i)} w_{i,j} \mathbf{x}_j \right\|^2$$

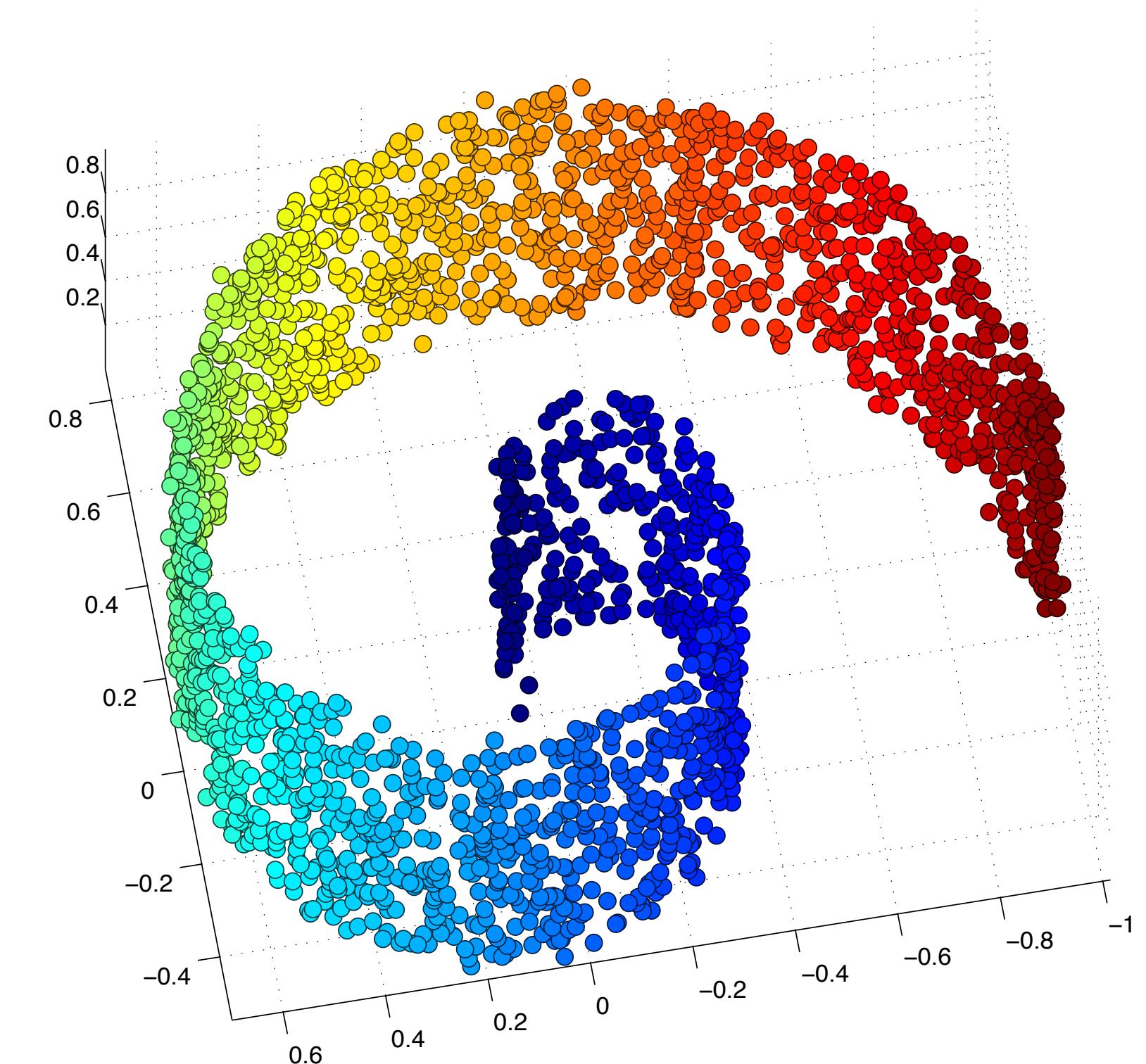


# Locally Linear Embedding

- The weights will work just as well in a lower dimensional space
  - We assume local linearity

$$C_L(\mathbf{W}) = \sum_i \left\| \mathbf{y}_i - \sum_{j \in N(i)} w_{i,j} \mathbf{y}_j \right\|^2$$

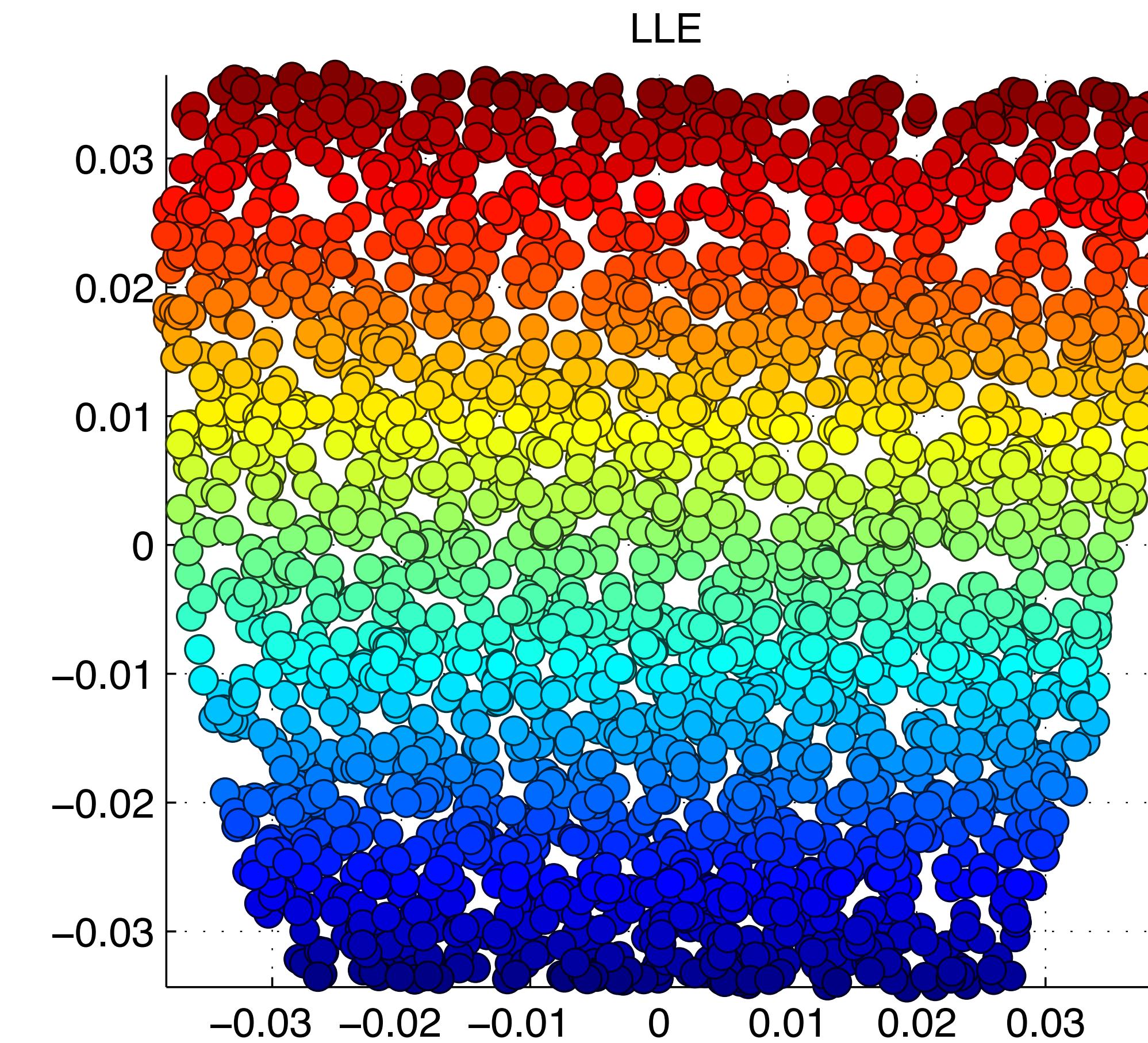
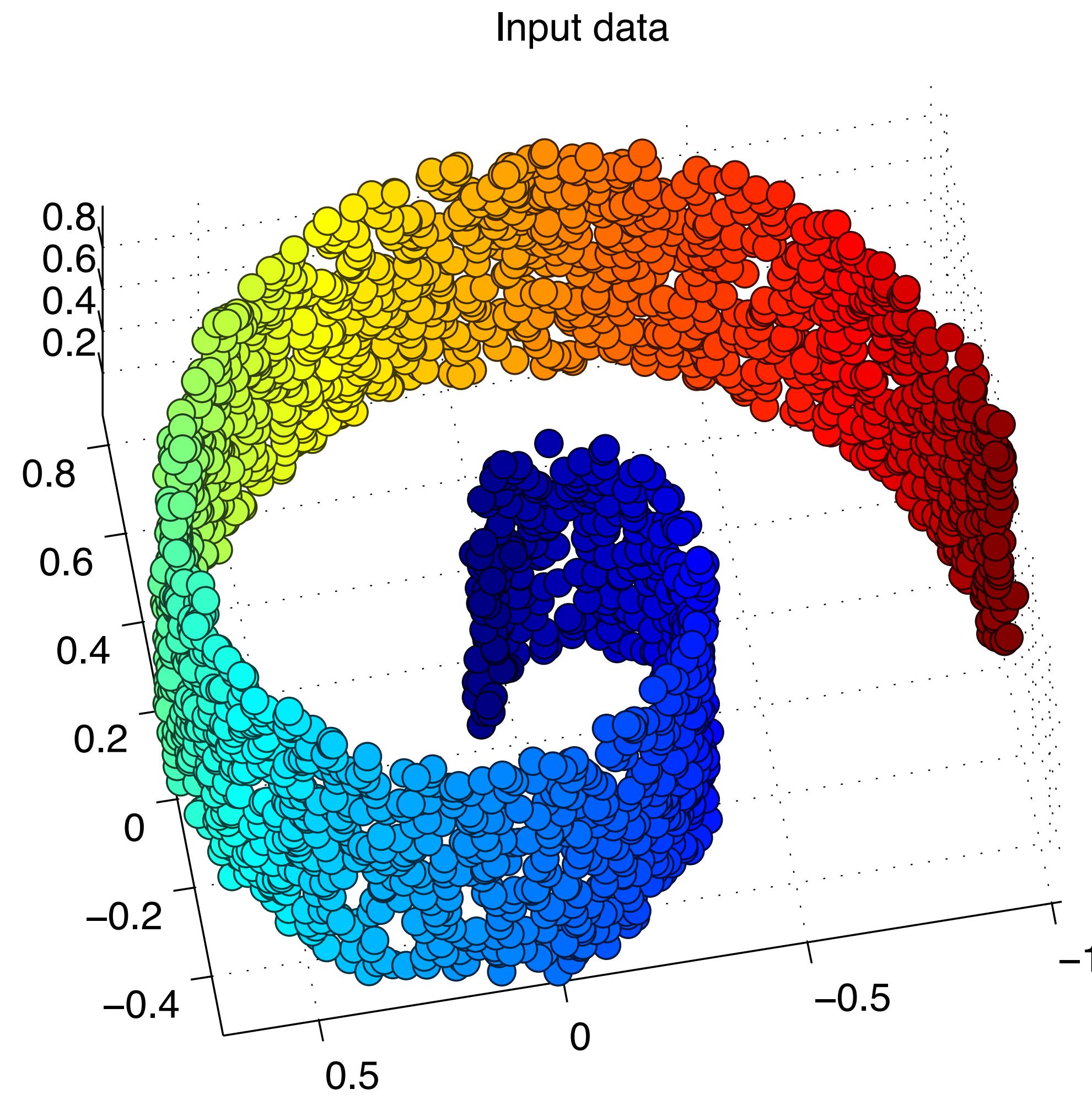
- We now need to find a  $\mathbf{y}$  that minimizes the above expression



# How to

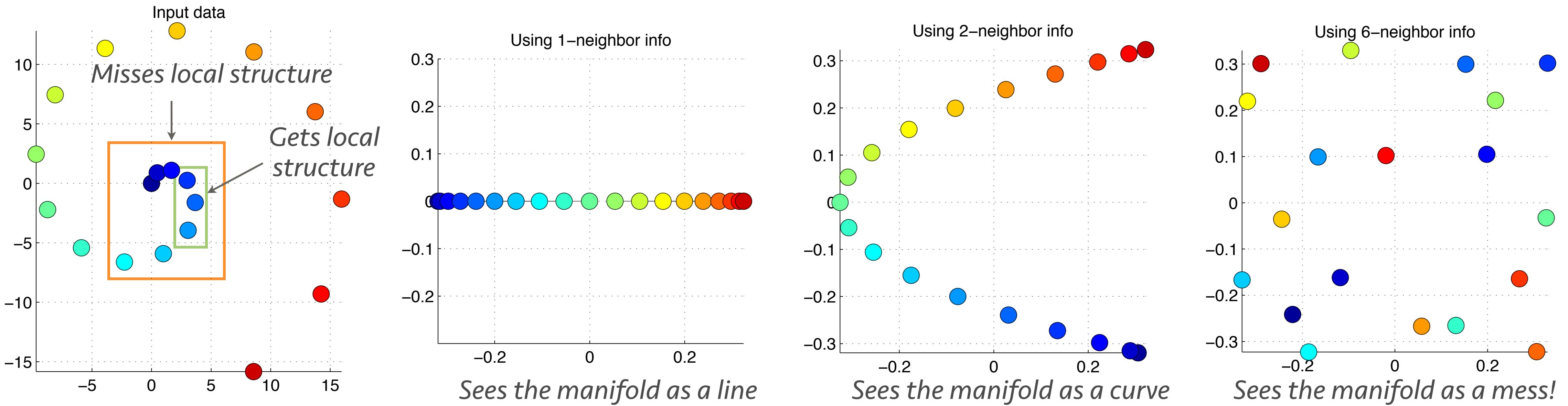
- Impose constraint on  $\mathbf{W}$ 
  - Rows must sum to 1
    - Make rotation, scale, translation invariant
- This becomes an eigendecomposition problem again!
  - This time we eigendecompose  $(\mathbf{I} - \mathbf{W})^\top \cdot (\mathbf{I} - \mathbf{W})$
  - Discard smallest eigenvector (has zero eigenvalue)
  - The rest of the smallest eigenvectors will contain variates of  $\mathbf{y}$

# Swiss roll example



# A note on manifold methods

- Try to highlight the local neighborhoods
  - e.g. clip the adjacency matrix, ignore distant points, etc.

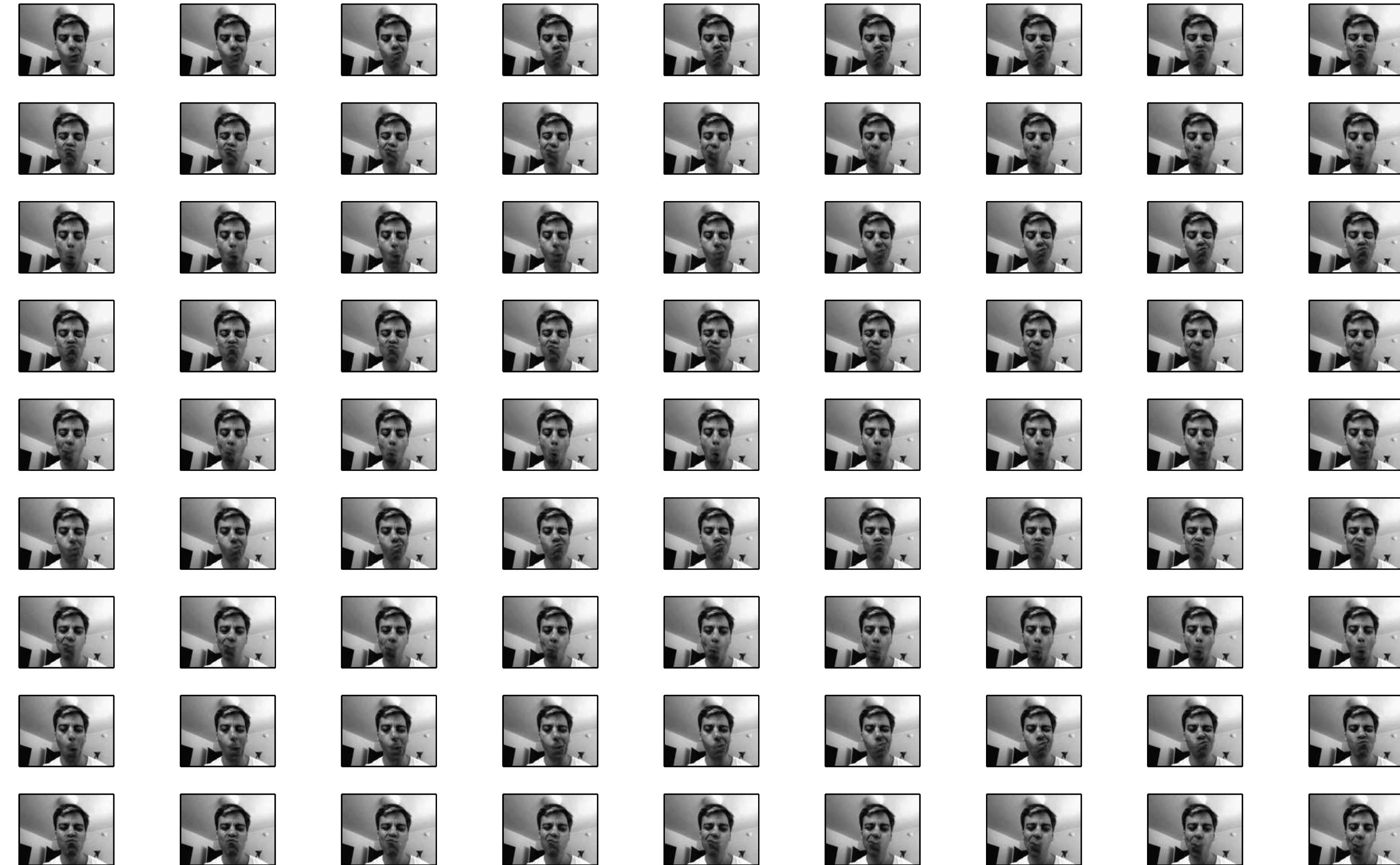


# A video example

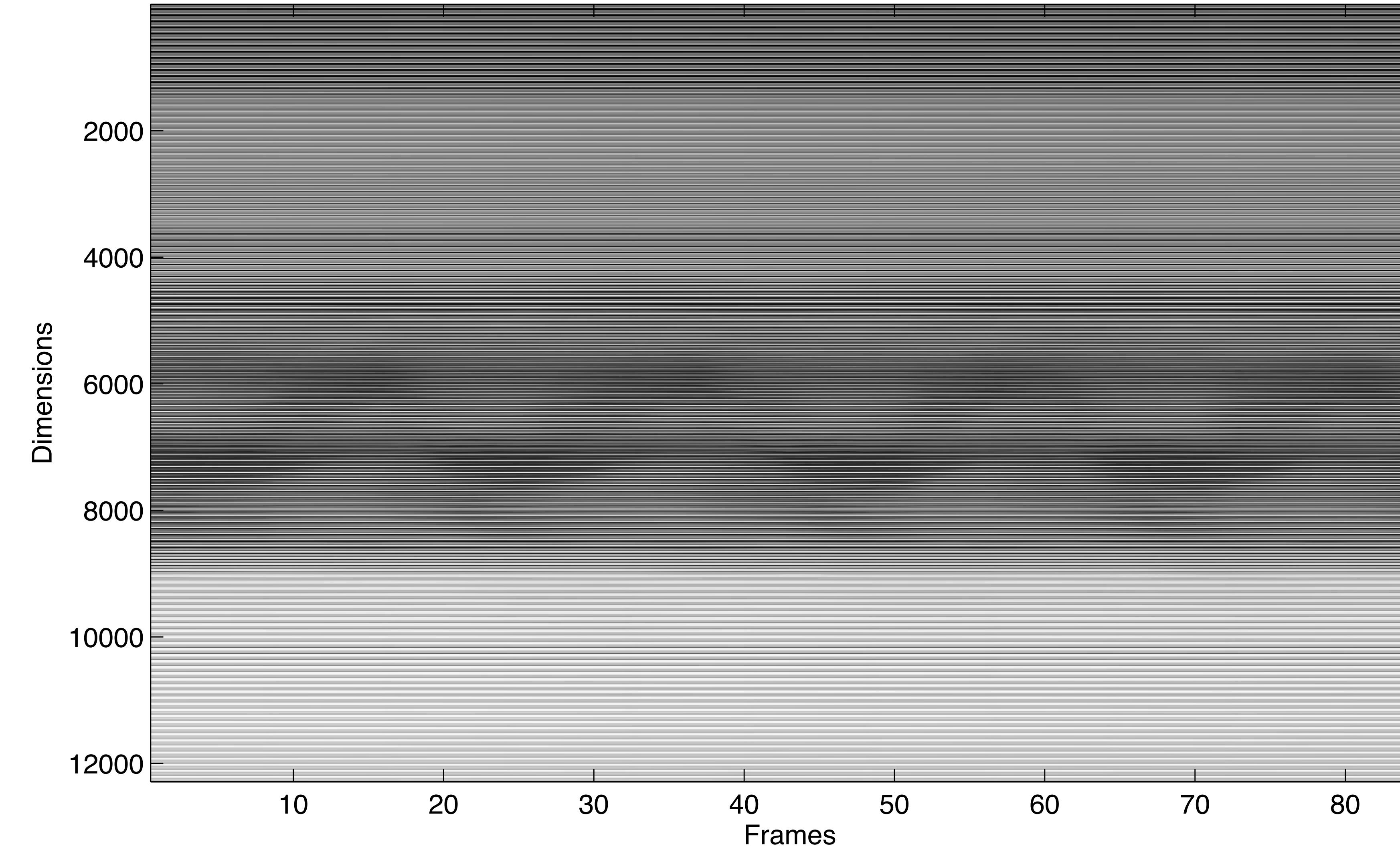
- A high dimensional input
  - $240 \times 320 \text{ pixels} = 76,800 \text{ dimensions}$
- Low dimensional structure
  - Moving lips around
- Can we simplify the data?



# Each frame is a pose



# When unraveled



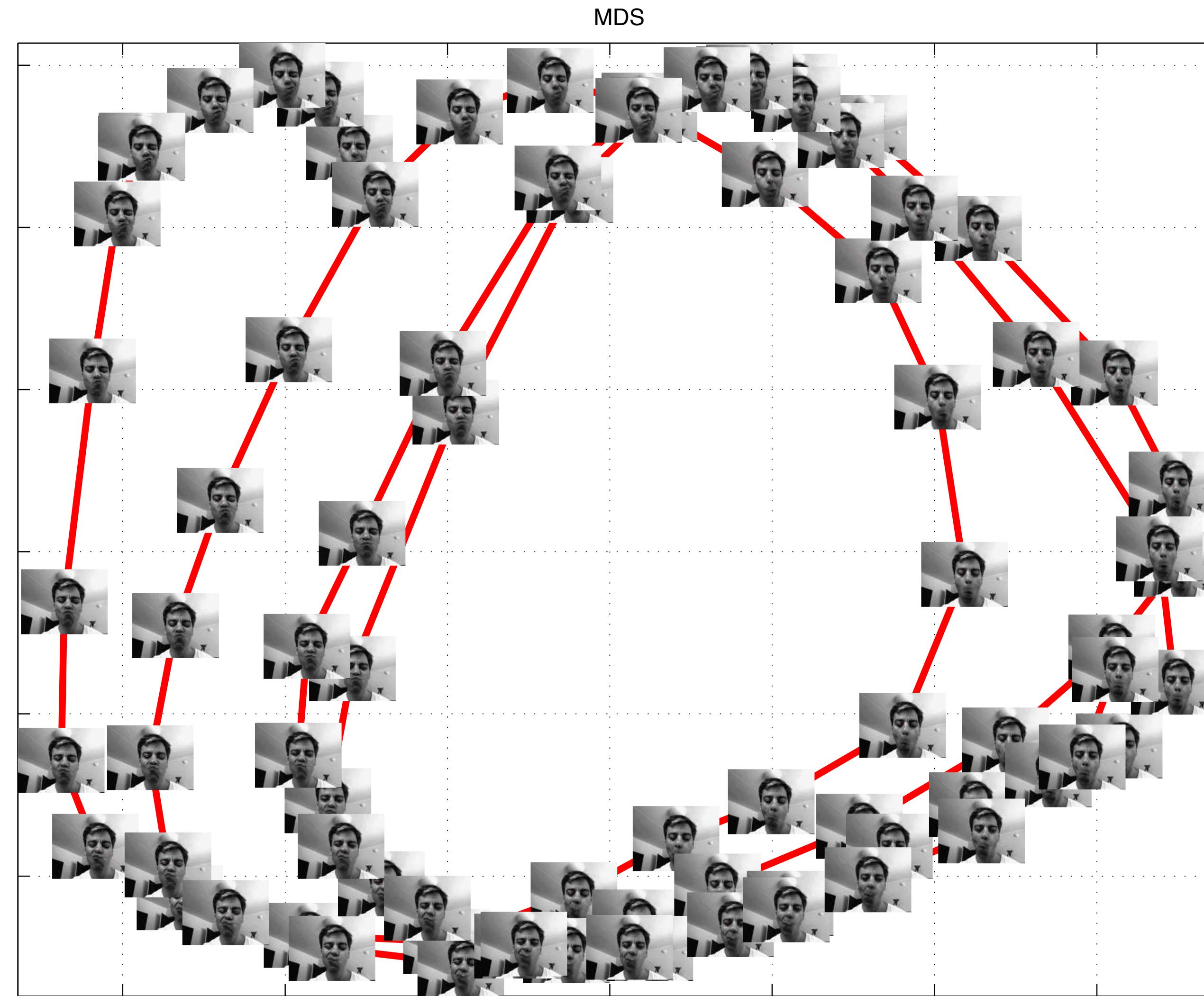
# The manifold

- We only have a couple of dimensions
  - The primary being the movements of the lips
- The actual space is much lower dimensional than the video is
  - It is smooth
  - And it is “circular”

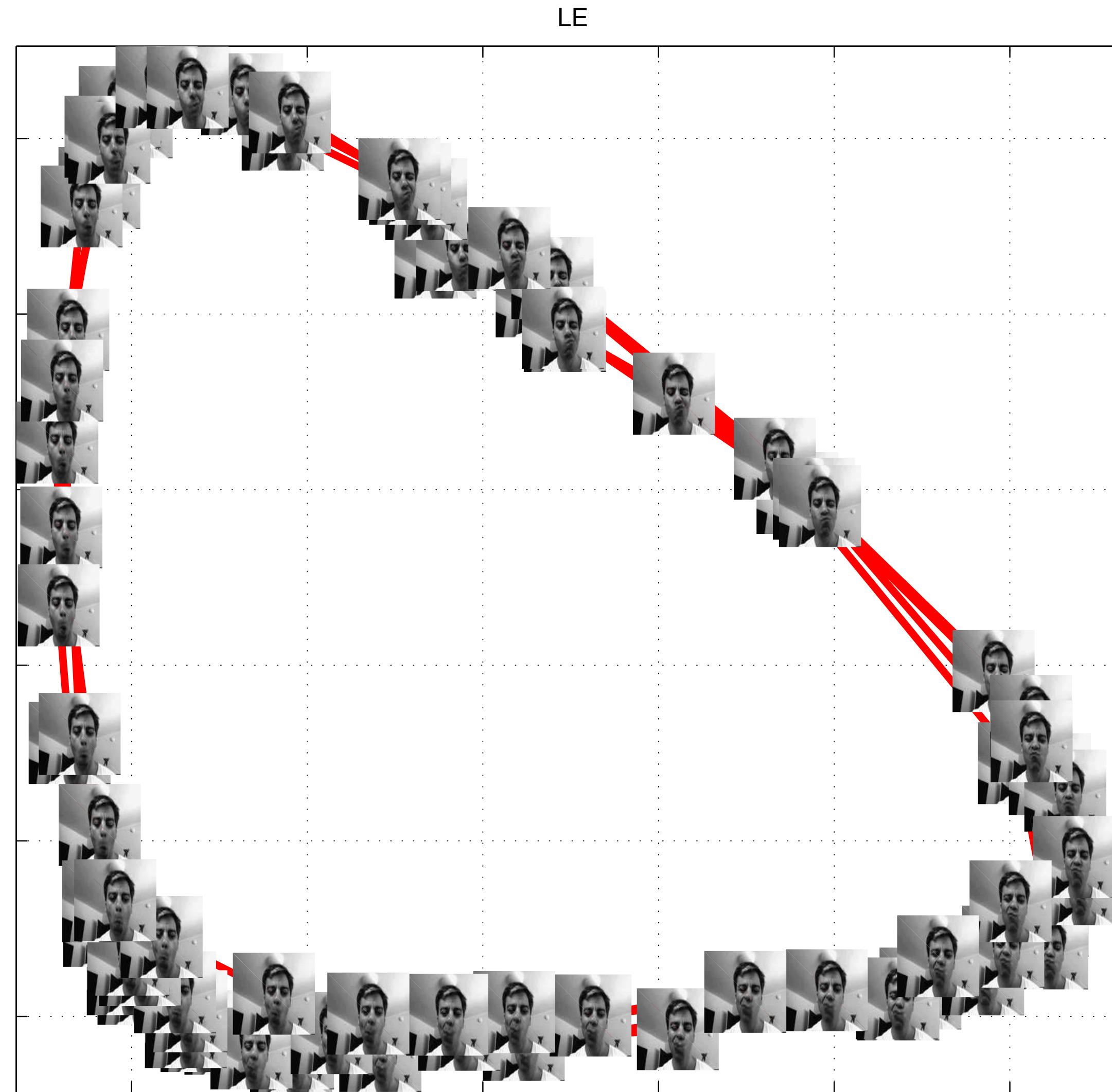
# PCA



# MDS

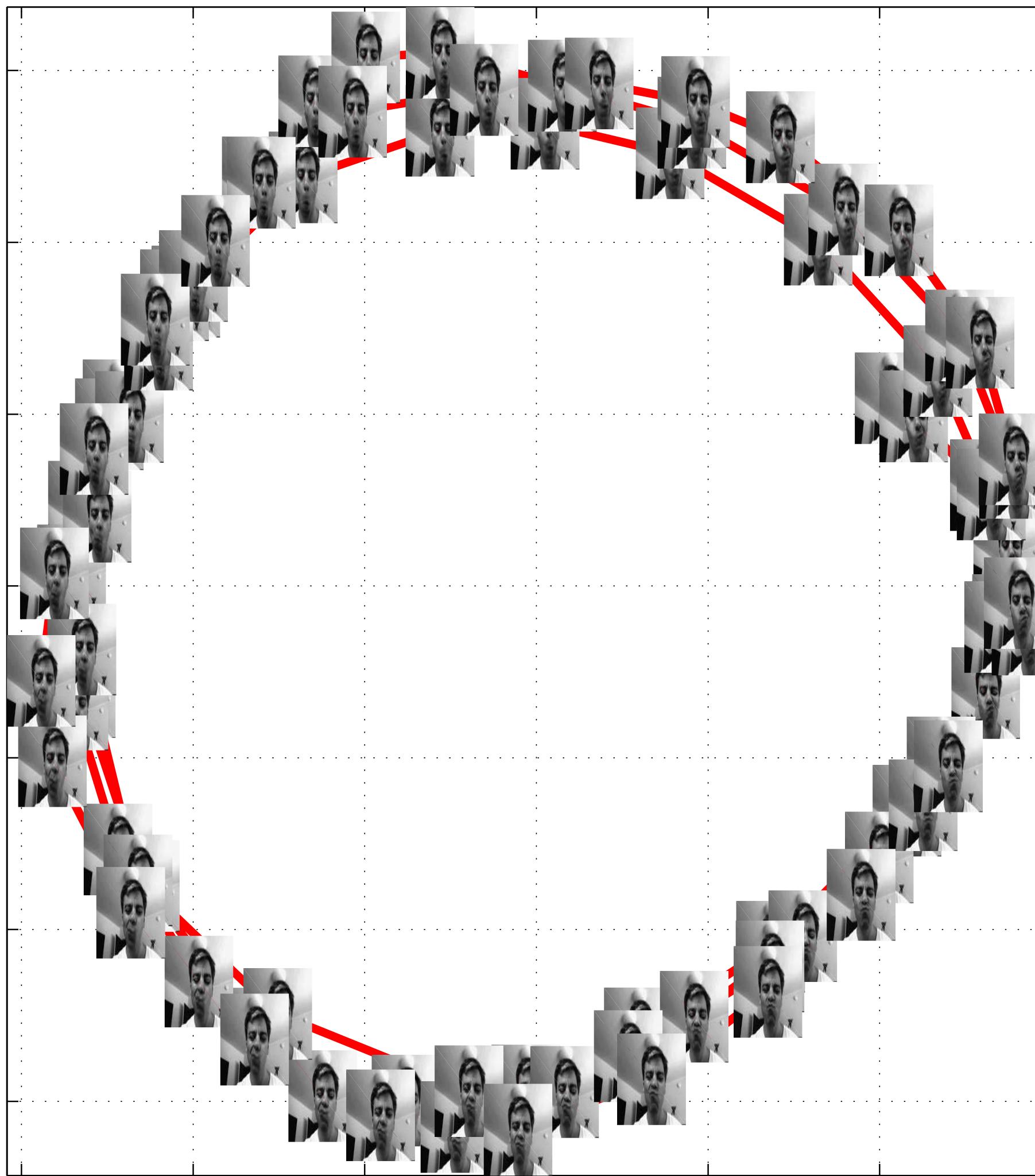


# Laplacian Eigenmaps

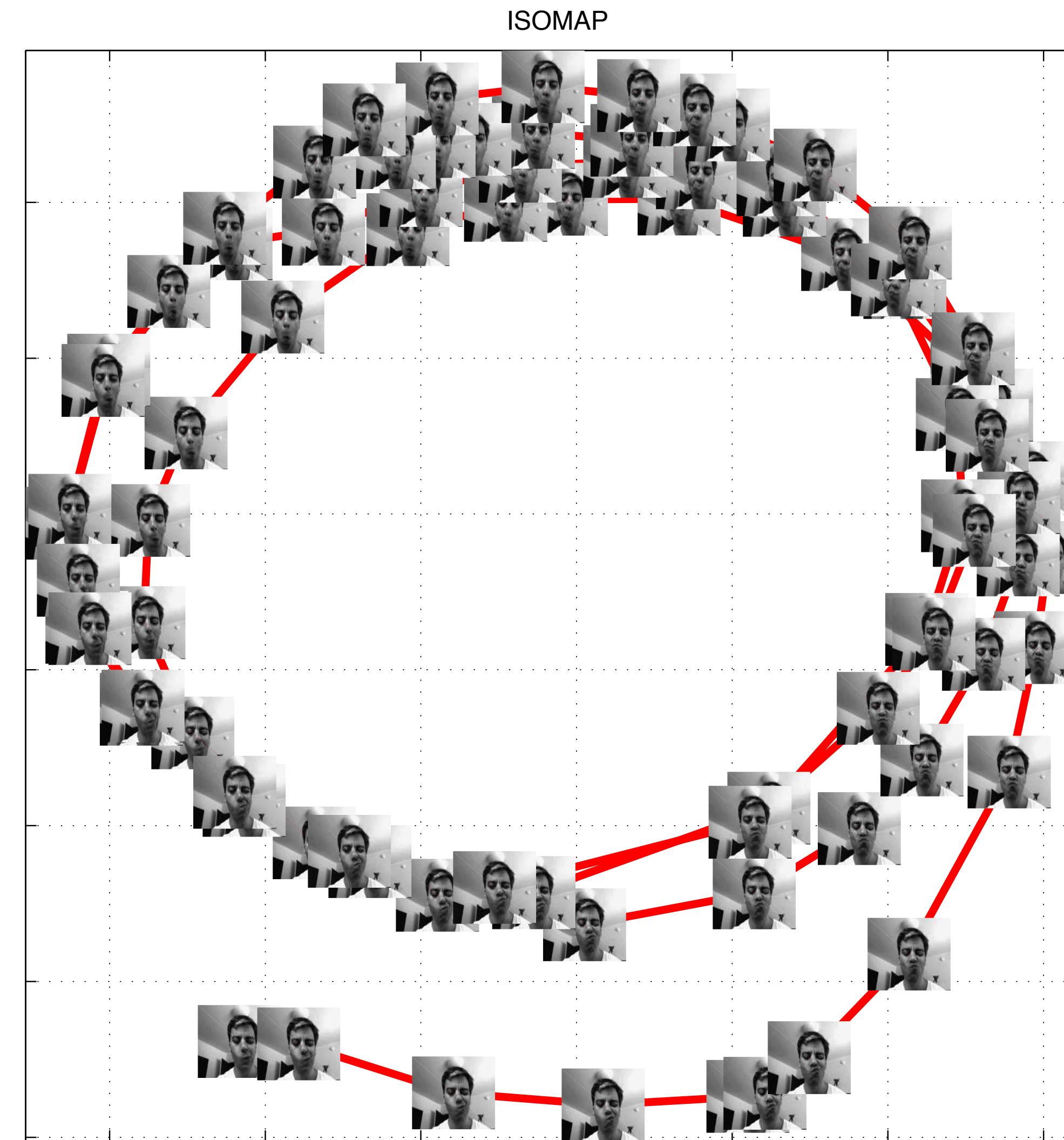


# LLE

LLE



# ISOMAP

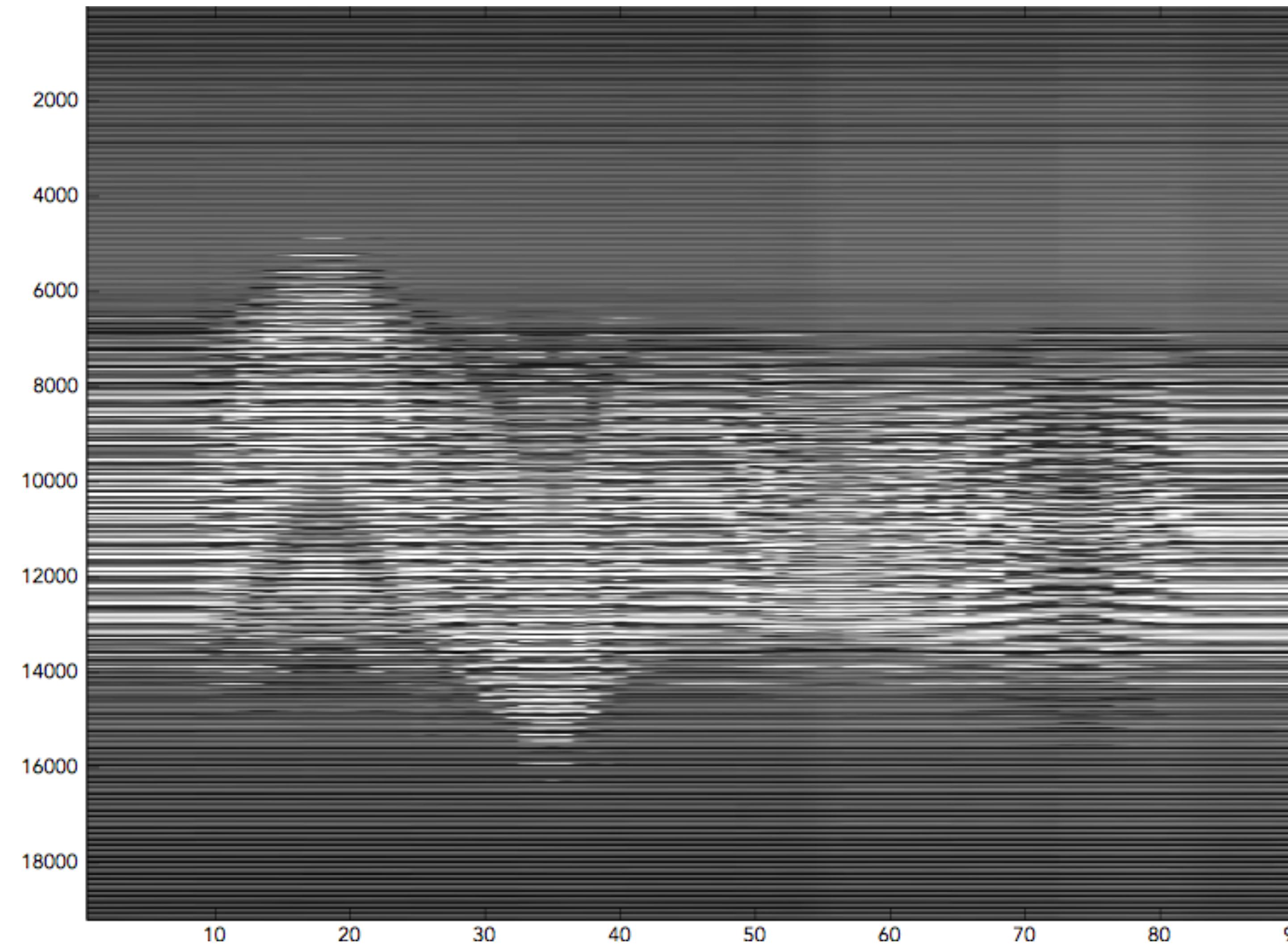


# Another example movie

- This time we have distinct axes
  - Up-Down-Left-Right head movement

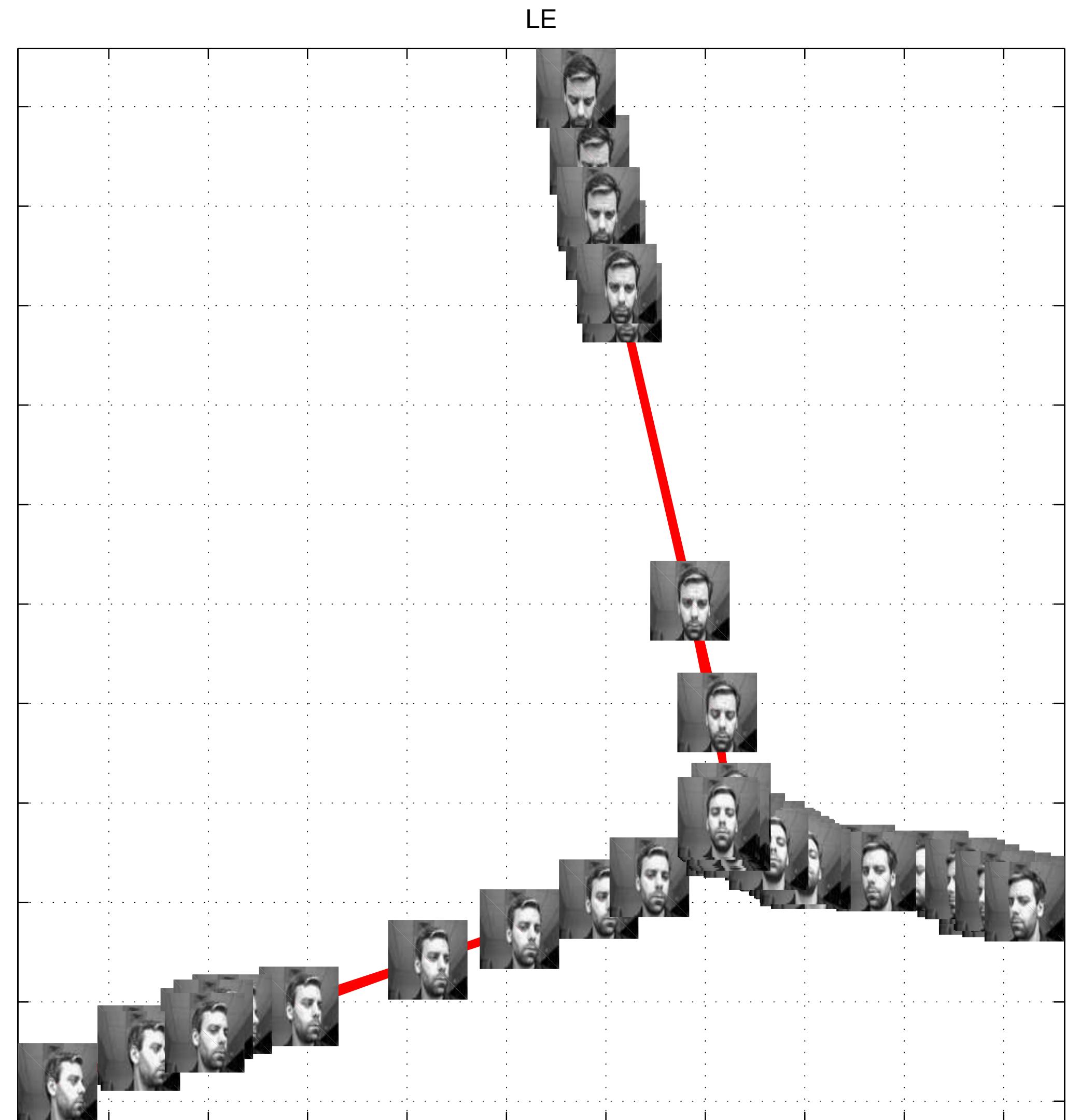


# Input as a matrix



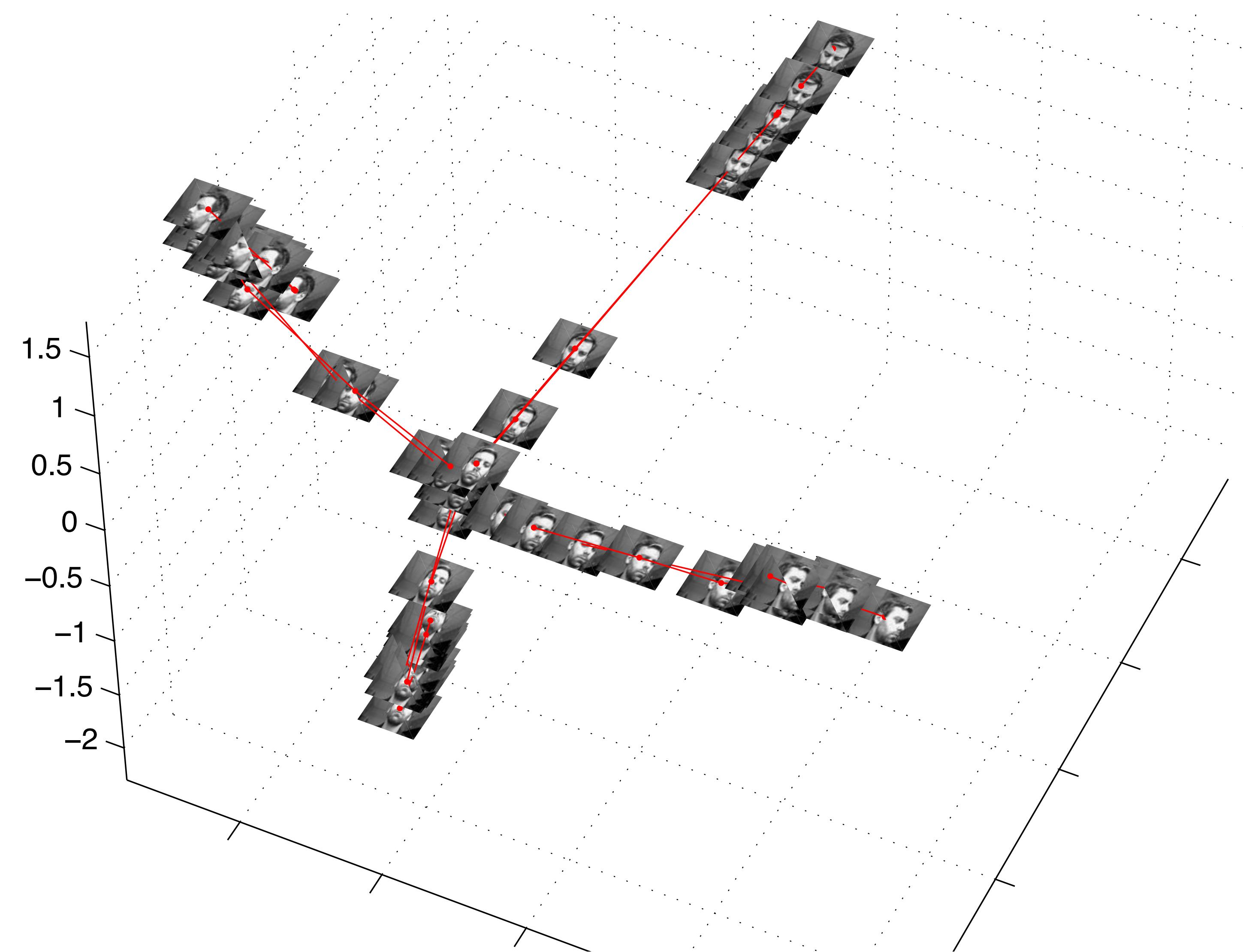
# Laplacian Eigenmaps 2D

- Stems coincide with head movements
- Missing dimension
  - Looking up



# Laplacian Eigenmaps 3D

- Better!



# The usual question

- Which method do I use?
  - The usual answer: “no good answer”
- There are tons of dimensionality reduction approaches, experimentation pays off!

# Recap

- Kernel PCA
  - Map to a non-linear space where things are more appropriate for PCA transforms
- MDS
  - Find an embedding that maintains distances
- Manifolds
  - ISOMAP
  - Locally Linear Embedding
  - Laplacian Eigenmaps

# Next lecture

- That's all for unsupervised learning for now!
- Moving to classification and supervised learning
- Starting from matched filtering and moving to more advanced approaches to classify data

# Reading

- Textbook section 6.7
- Kernel PCA (optional)
  - <http://www.springerlink.com/content/w0t1756772h41872/>
- Manifolds (optional)
  - <http://science.sciencemag.org/content/290/5500/2268>
  - <http://www.sciencemag.org/cgi/reprint/290/5500/2323.pdf>
  - <http://www.sciencemag.org/cgi/reprint/290/5500/2319.pdf>