

CS49000-VIZ - Fall 2020

Introduction to Data Visualization

Dimensionality

Reduction

Lecture 16

November 12, 2020

Outline

High-dimensional data

Dimensionality reduction

Manifold learning

Computational aspects

Why High-Dimensional?

- Data samples with many attributes
 - Tables with many columns, medical records, ...
- Amalgamate individual properties into high-dimensional *feature vectors*

Challenges

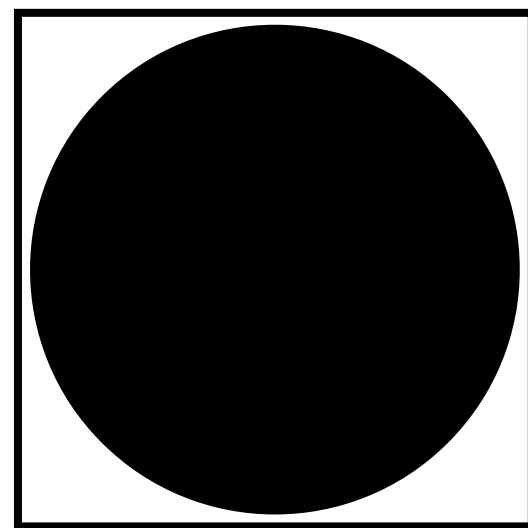
- “Curse of dimensionality”
 - Sampling becomes exponentially costly
 - Space accumulates in “corners” of hypercube
 - Data processing becomes extremely expensive / intractable
- How to visualize?
 - Missing intuition

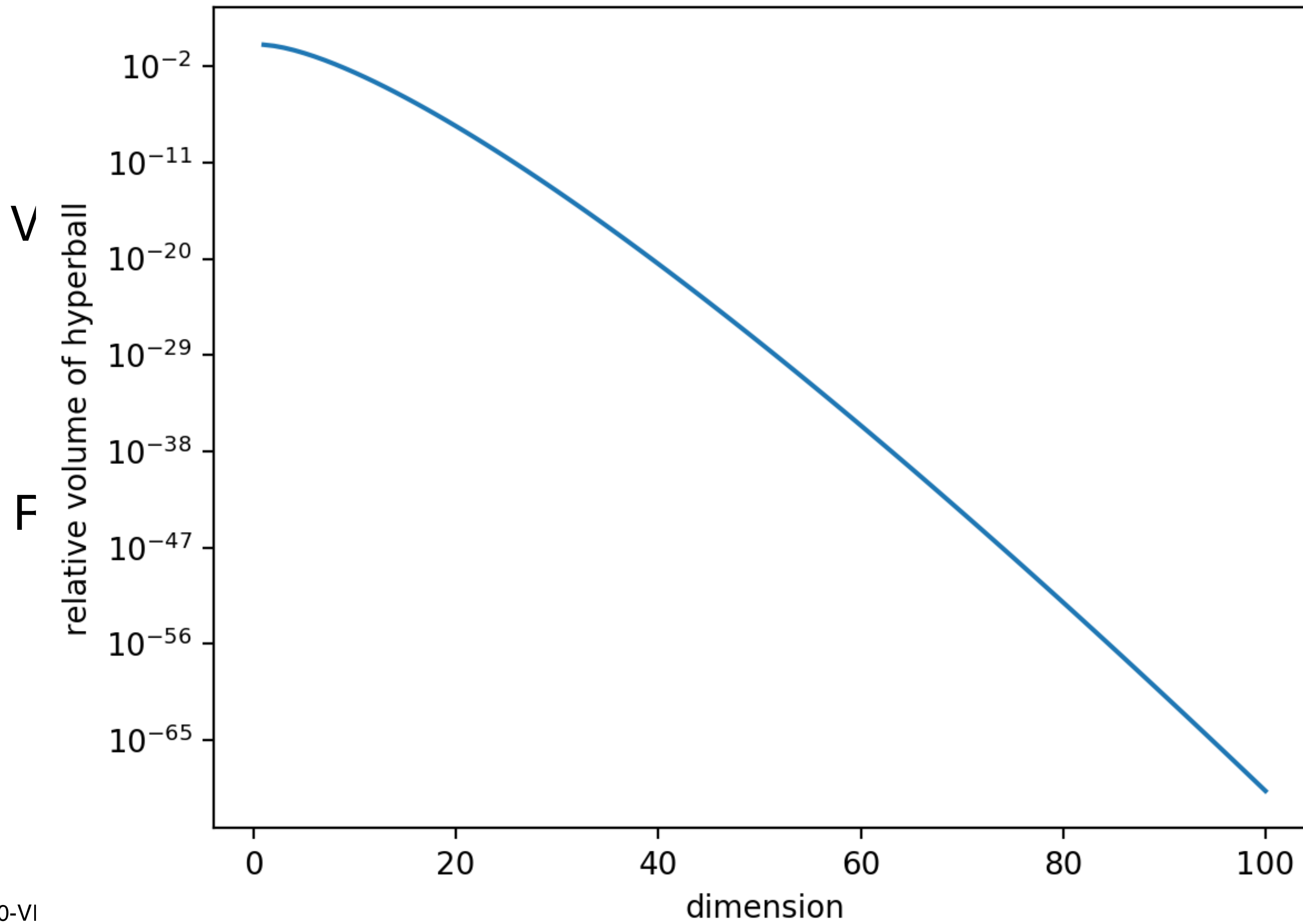
Curse of Dimensionality

Volume of unit hyper-ball in n-dimensions: $V_n = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$

$$V_{2k} = \frac{\pi^k}{k!} \quad V_{2k+1} = \frac{2(k!)(4\pi)^k}{(2k+1)!}$$

Ratio to volume of bounding hyperbox $\frac{V_{2k}}{2^k} = \left(\frac{\pi}{4}\right)^k \frac{1}{k!}$





Dimension Reduction

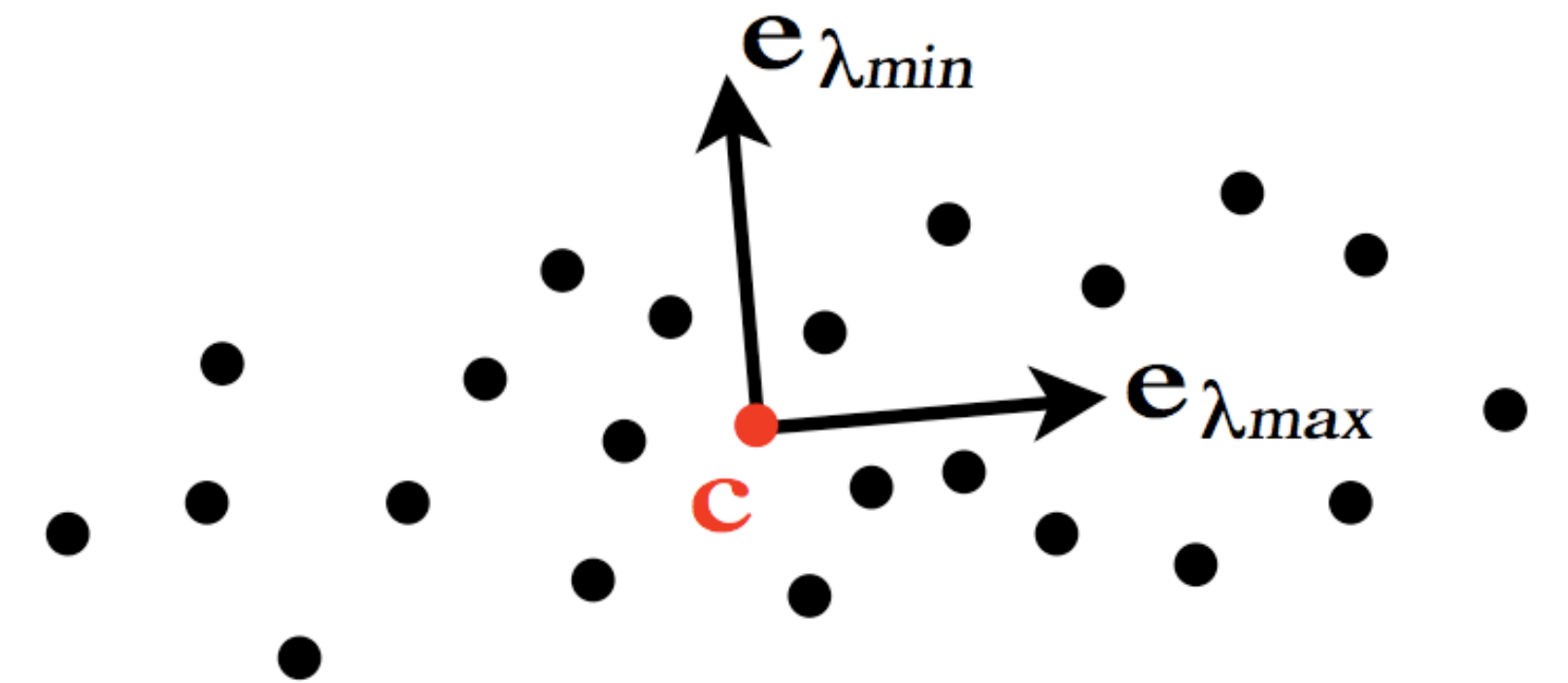
PCA

- Assume a point cloud $(\mathbf{x}_i)_{i=1,\dots,n} \in \mathbb{R}^k$
- Interpret these points as observations of a random variable $\mathbf{X} \in \mathbb{R}^k$
- The empirical mean (centroid) is $\mathbf{c} = \frac{1}{n} \sum_i \mathbf{x}_i$
- The covariance matrix is given by $\mathbf{A} = \bar{\mathbf{X}}\bar{\mathbf{X}}^T$

$$A_{jl} = \frac{1}{n} \sum_i (x_{ij} - c_j)(x_{il} - c_l)$$

PCA

- The eigenvectors of the covariance matrix form a data-dependent coordinate system
- The first m eigenvectors (in decreasing order of associated eigenvalues) span the m principal dimensions of the point cloud.



MDS

Multidimensional Scaling

- Input: dissimilarity matrix $\Delta \in \mathbb{R}^{n \times n}$

$$\Delta_{ij} = \delta_{i,j} \quad \text{with } \Delta_{ij} > 0 \text{ and } \Delta_{ii} = 0$$

- Goal: find $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ such that

$$\forall (i, j) \quad \|\mathbf{x}_i - \mathbf{x}_j\| \approx \delta_{i,j}$$

MDS

Method:

- Center(*): $B = -\frac{1}{2}H\Delta H$ where $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$
- Spectral decomposition: $B = U\Lambda U^T$
- Clamp(*): $(\Lambda_+)_{ij} = \max(\Lambda_{ij}, 0)$
- Solution: first d columns of $X = U\Lambda_+^{\frac{1}{2}}$

This is a PCA!

(*): if Δ measures Euclidean distances, B is positive semidefinite

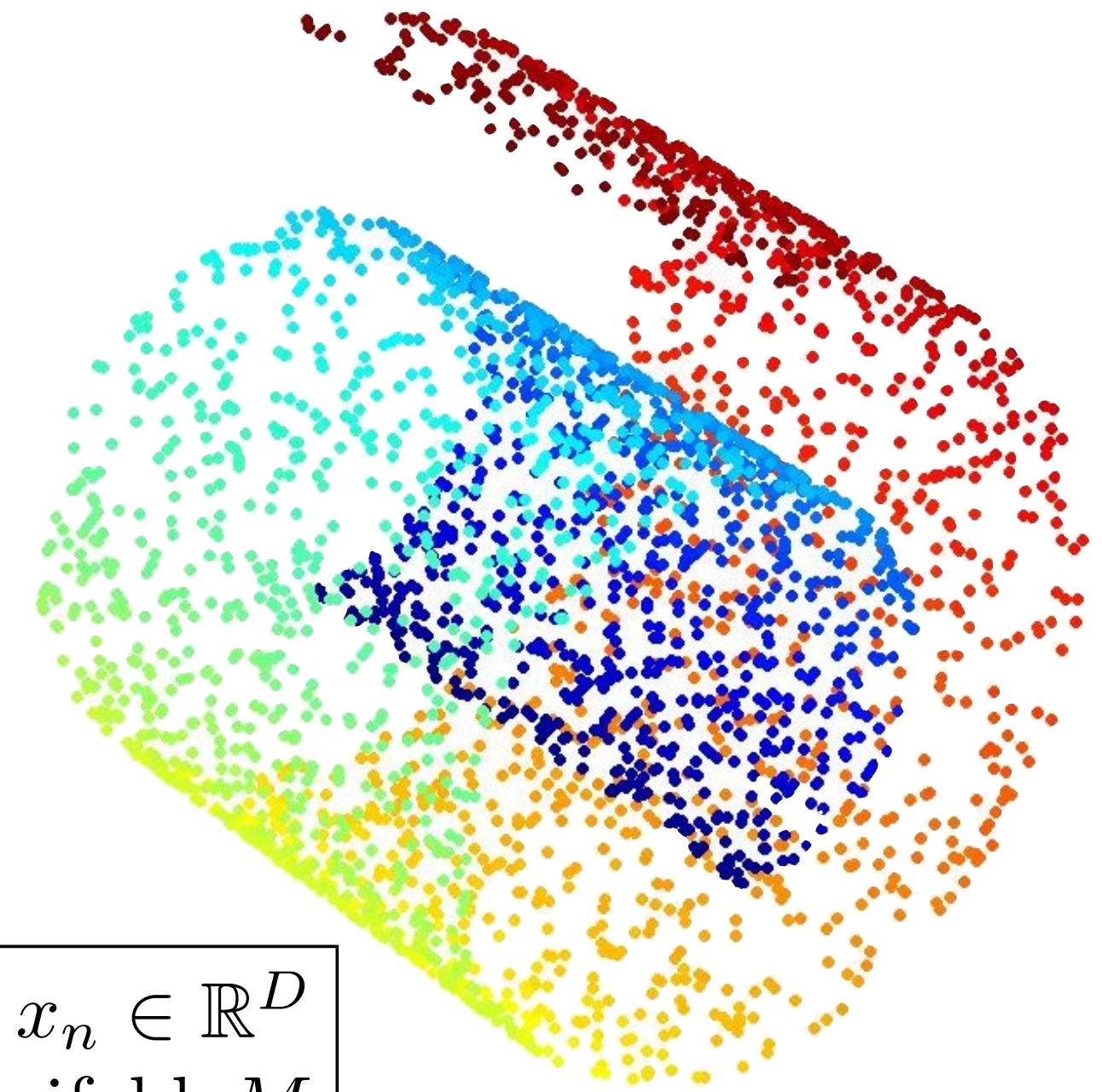
Manifold Learning

- Move away from linear assumption
- Assume curved low-dimensional geometry
- Assume smoothness

➔ MANIFOLD

- Manifold learning

Problem: Given points $x_1, \dots, x_n \in \mathbb{R}^D$ that lie on a d -dimensional manifold M that can be described by a single coordinate chart $f : M \rightarrow \mathbb{R}^d$, find $y_1, \dots, y_n \in \mathbb{R}^d$, where $y_i \stackrel{\text{def}}{=} f(x_i)$.



ISOMAP

Tenenbaum et al., Science, 2000

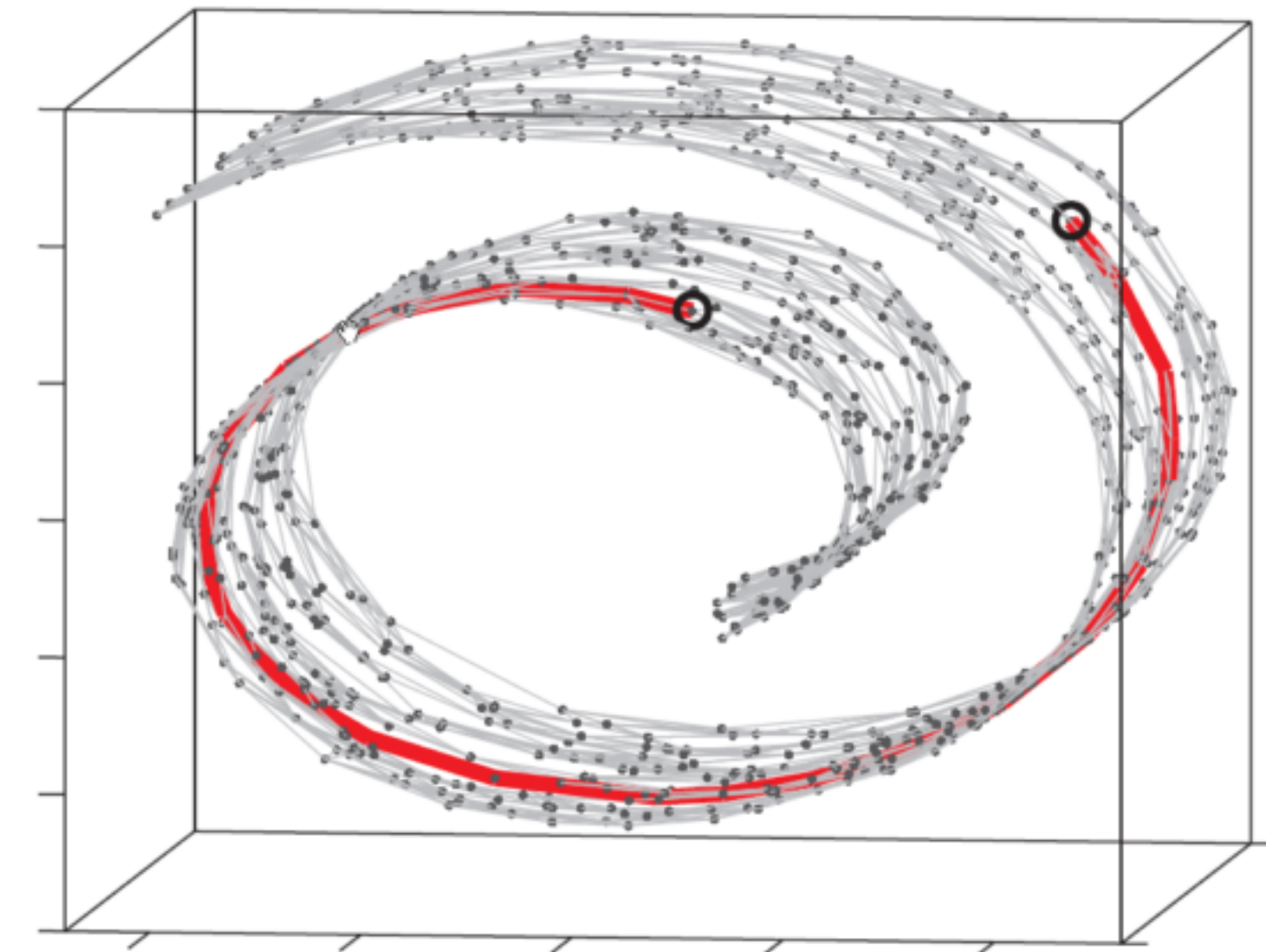
- Form k-NN graph of input points
- Form dissimilarity matrix Δ as square of approximated geodesic distance between points
- Compute MDS!

ISOMAP

- Computing geodesic distance
 - standard graph problem in CS

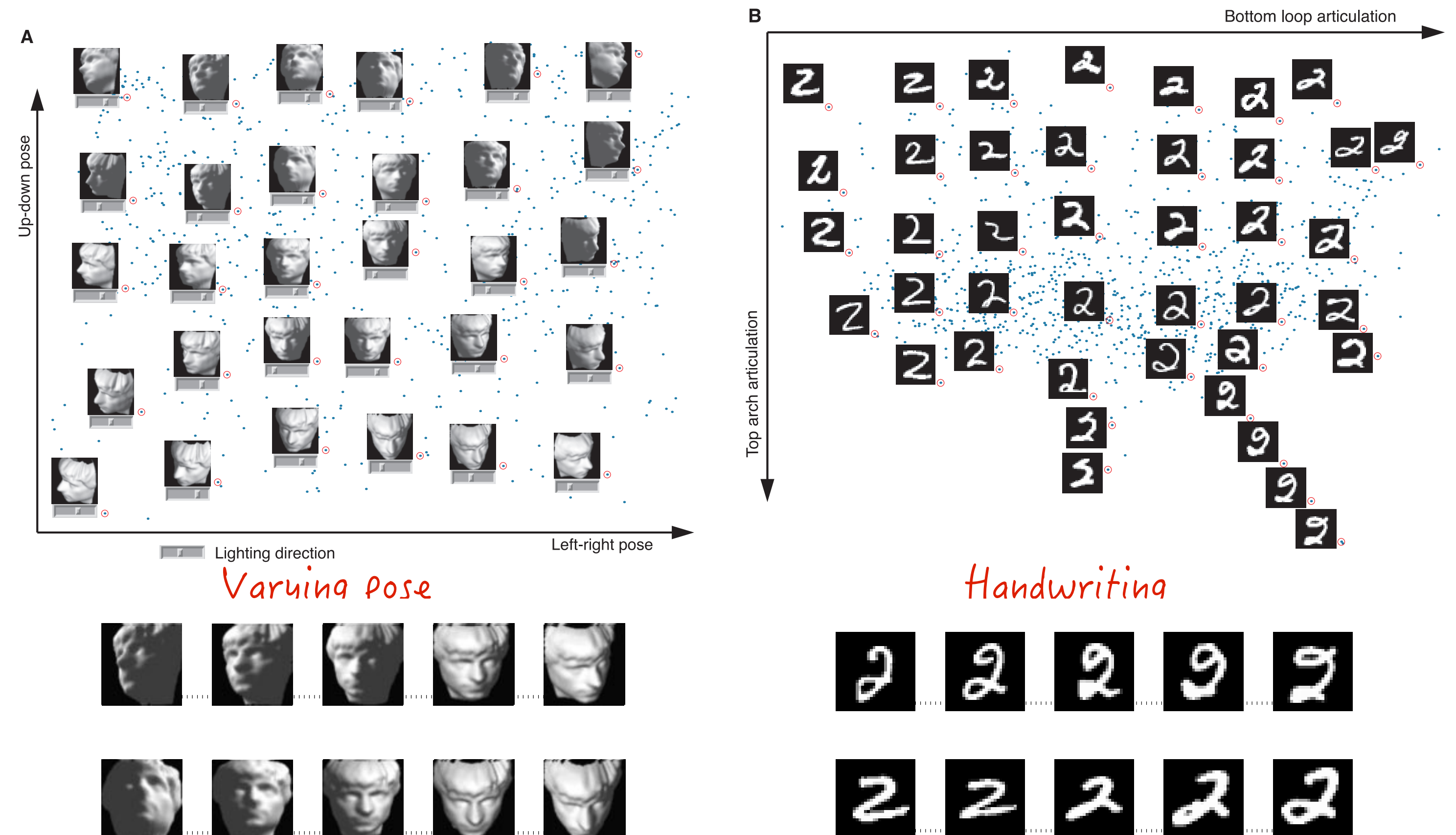
Dijkstra Algorithm $O(|E| + |V| \log |V|)$

```
1 function Dijkstra(Graph, source):
2   for each vertex v in Graph:           // Initializations
3     dist[v] := infinity                 // Unknown distance function from source to v
4     previous[v] := undefined           // Previous node in optimal path from source
5   dist[source] := 0                    // Distance from source to source
6   Q := the set of all nodes in Graph   // All nodes in the graph are unoptimized - thus are in Q
7   while Q is not empty:               // The main loop
8     u := vertex in Q with smallest dist[]
9     remove u from Q
10    for each neighbor v of u:          // where v has not yet been removed from Q.
11      alt := dist[u] + dist_between(u, v) // be careful in 1st step - dist[u] is infinity yet
12      if alt < dist[v]:                 // Relax (u,v,a)
13        dist[v] := alt
14        previous[v] := u
15  return previous[]
```



ISOMAP

- Successful in computer vision problems



Locally Linear Embedding

Roweis et al., Science, 2000

Intuition

- smooth manifold is locally close to linear

Method

- Characterize local linearity through weight matrix W such that $\|\mathbf{x}_i - \sum_{j \in \mathcal{N}(i)} W_{ij} \mathbf{x}_j\|^2$ is minimized
- Find d -dimensional \mathbf{y} 's that minimize $\sum_i \|\mathbf{y}_i - \sum_j W_{ij} \mathbf{y}_j\|^2$
- Solution is obtained as first d eigenvectors of $(I - W)^T (I - W)$

Locally Linear Embedding

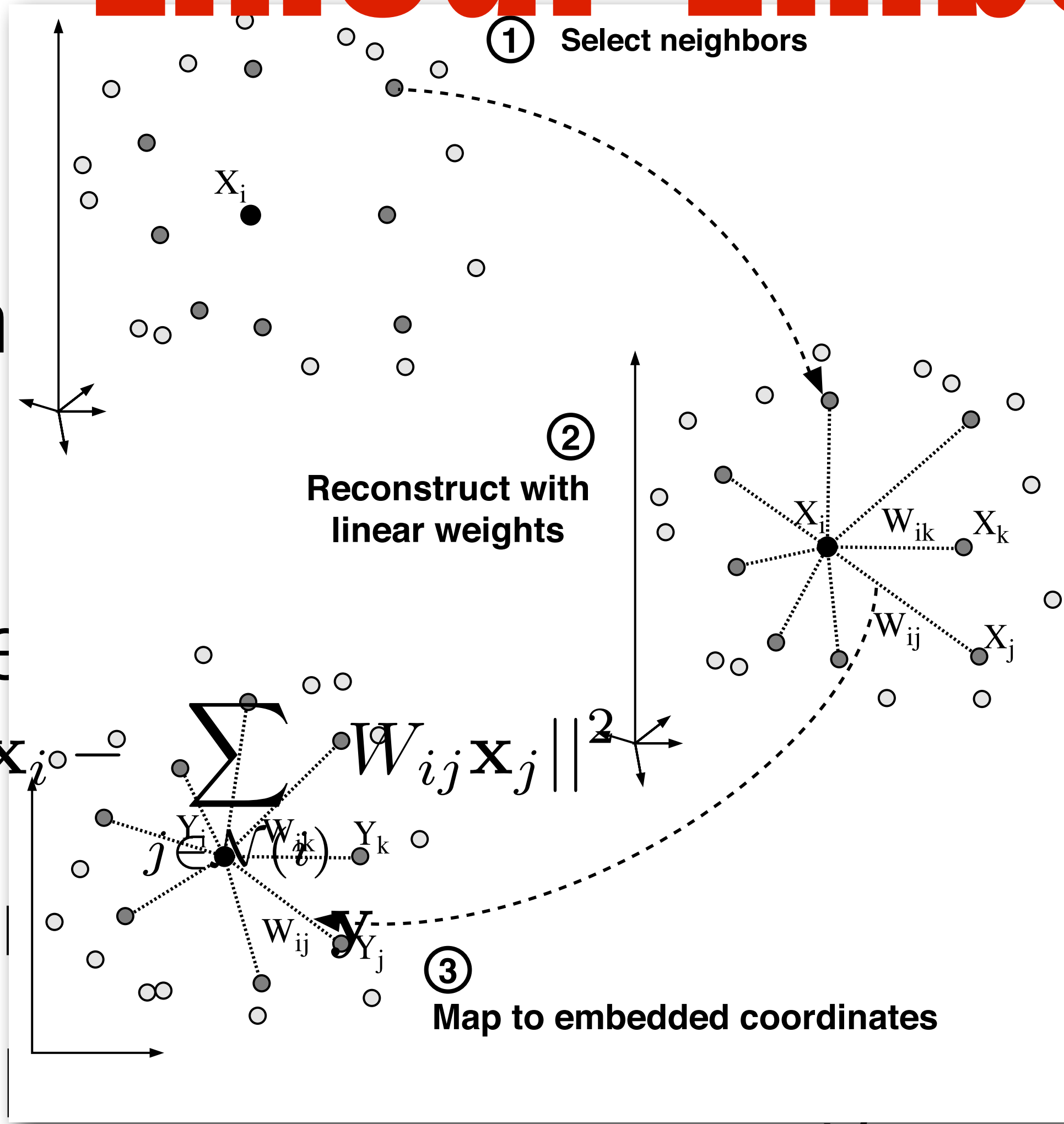
Roweis et al., Science, 2000

Intuition

- smooth manifold

Method

- Characterize such that $\|x_i - \sum_{j \in N(i)} W_{ij} x_j\|^2$
- Find d-dimensional
- Solution is on



near

weight matrix W
 sized

$$\sum_i \|y_i - \sum_j W_{ij} y_j\|^2$$

factors of

$$(I - W)^T (I - W)$$

Locally Linear Embedding

Roweis et al., Science, 2000

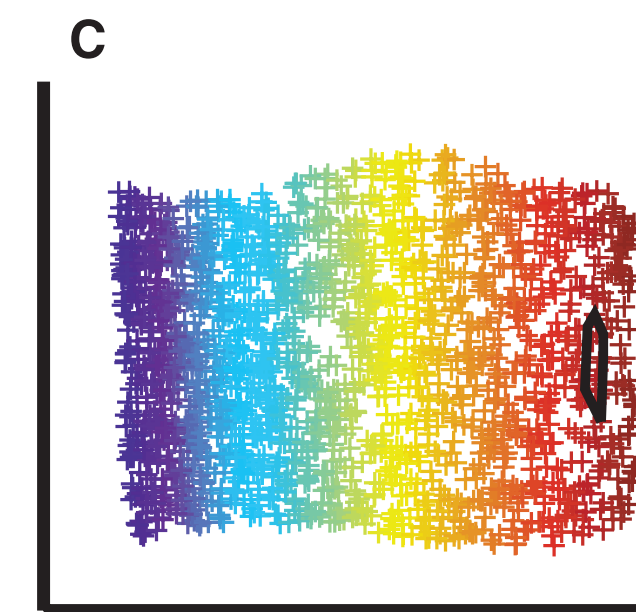
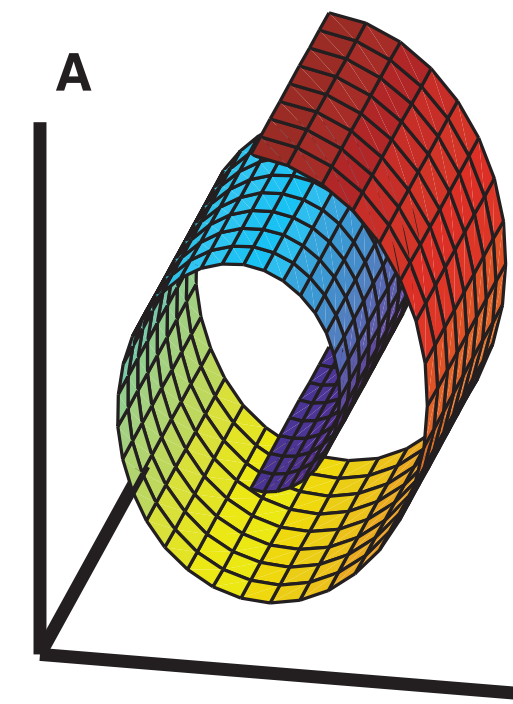
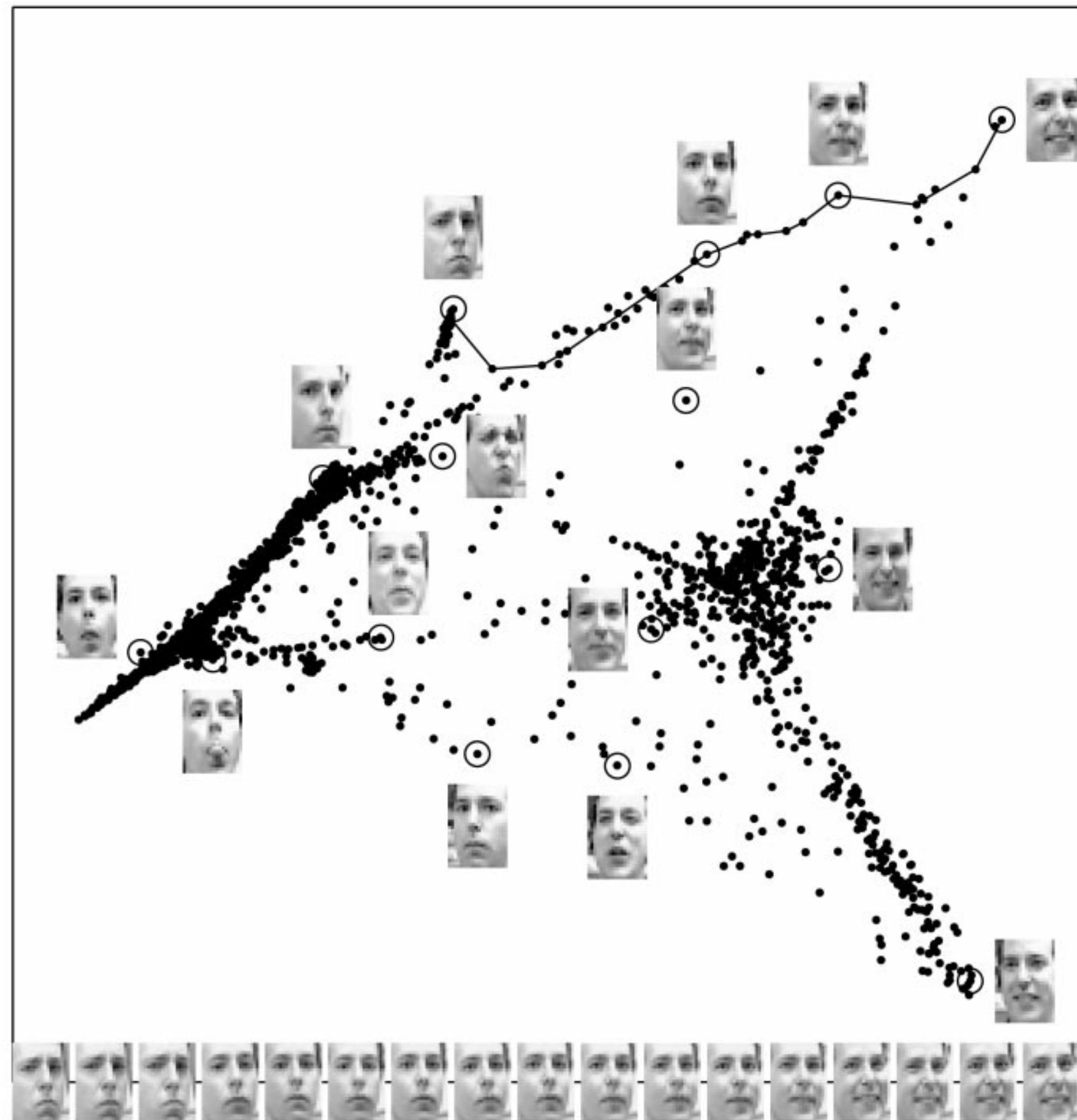
Intuition

- smooth manifold is locally close to linear

Method

- Characterize local linearity through weight matrix W such that $\|\mathbf{x}_i - \sum_{j \in \mathcal{N}(i)} W_{ij} \mathbf{x}_j\|^2$ is minimized
- Find d -dimensional \mathbf{y} 's that minimize $\sum_i \|\mathbf{y}_i - \sum_j W_{ij} \mathbf{y}_j\|^2$
- Solution is obtained as first d eigenvectors of $(I - W)^T (I - W)$

LLE



Laplacian Eigenmaps

Belkin, Niyogi, Advances in Neural Information Processing Systems, 2001

- Graph Laplacian of W is $L = D - W$

with D diagonal and $D_{ii} = \sum_j W_{ij}$

- Define W as either binary indicator of connectivity or through heat kernel
- Solve for $\sum_{i,j} W_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|^2 = \text{tr}(Y^T LY)$
- Compute eigenvectors \mathbf{y} of Laplacian associated with non-zero eigenvalue

Laplacian Eigenmaps

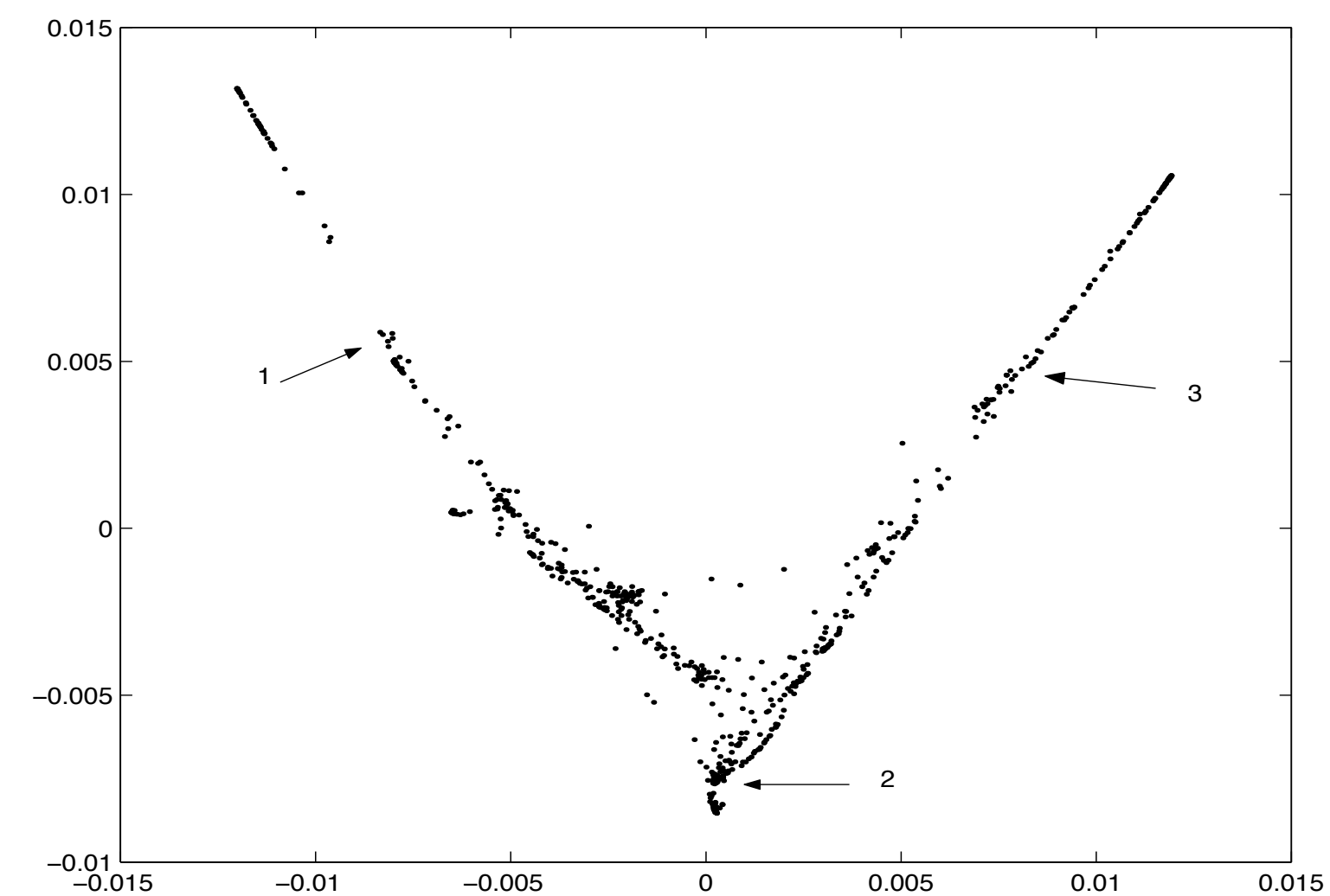
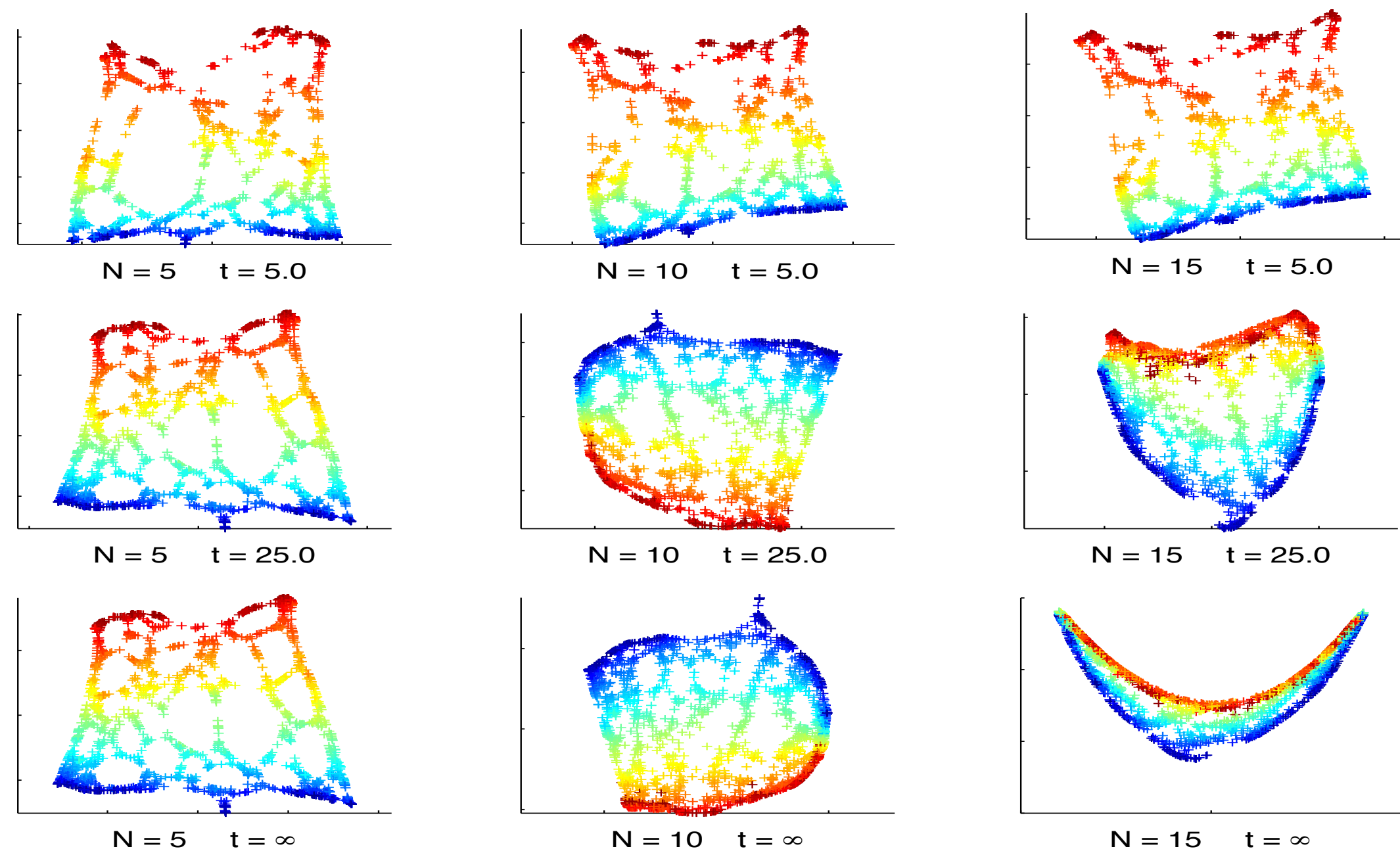
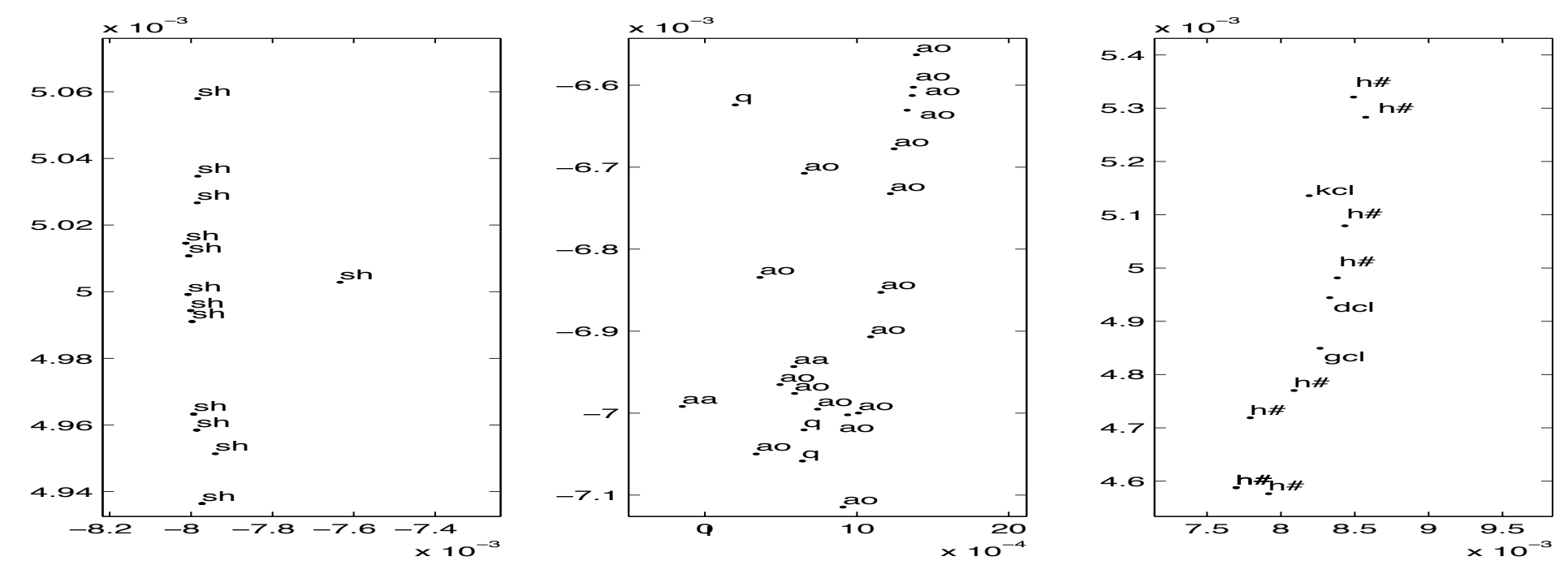


Figure 2: Two-dimensional representations of the “swiss roll” data, for different values of the number of nearest neighbors N and the heat kernel parameter t . $t = \infty$ corresponds to the discrete weights.



Computation

- Spectral decomposition intractable for large problems
- Approximate method can be used
- Nystrom method + column sampling