

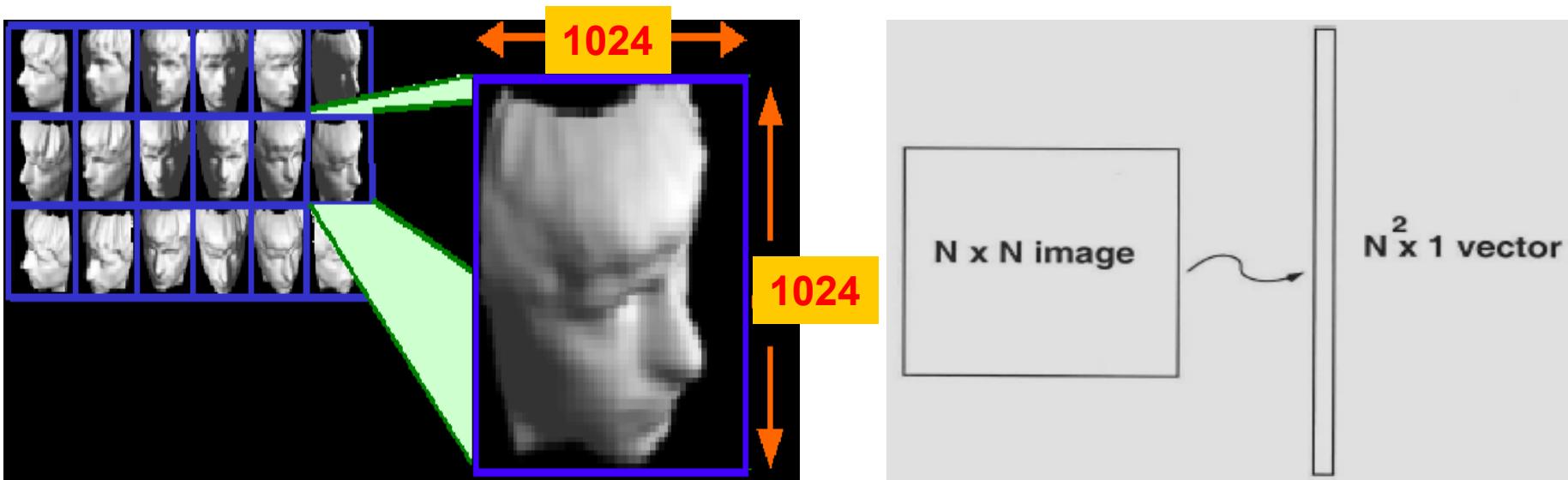
Dimensionality Reduction

Evgeny Burnaev
Skoltech

Dimensionality Reduction Problem

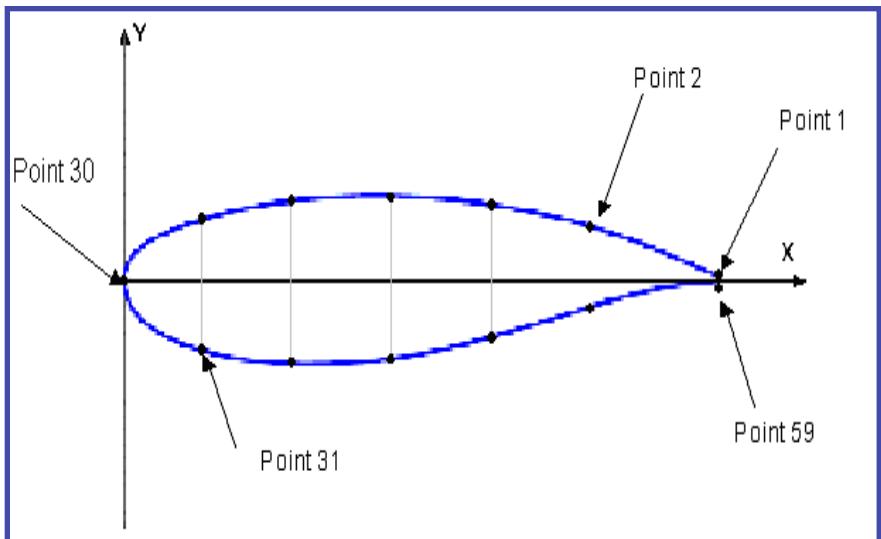
- object **O** is described by a **p**-dimensional vector $\mathbf{X(O)} \in \mathbb{R}^p$. Components of $\mathbf{X(O)}$ are features of **O**

Example 1 (human face): grey-scale pixel-wise representation



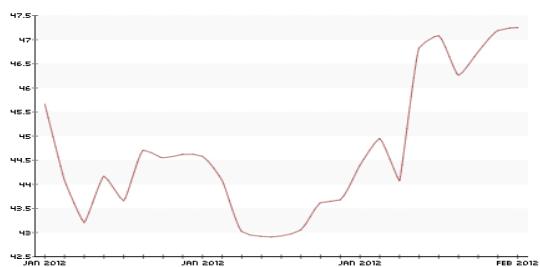
Face is represented by 1024×1024 pixels – dimension $p = 10^{20} \sim 1\ 000\ 000$

Example 2 (wing airfoil description).

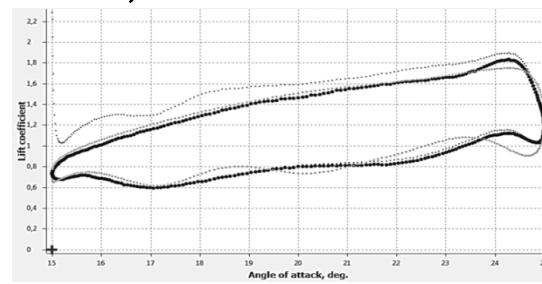


Airfoil $O \rightarrow X(O) = (x_1, x_2, \dots, x_p)^T$ - ordinates of upper and lower contours, $p \sim 50 \div 200$

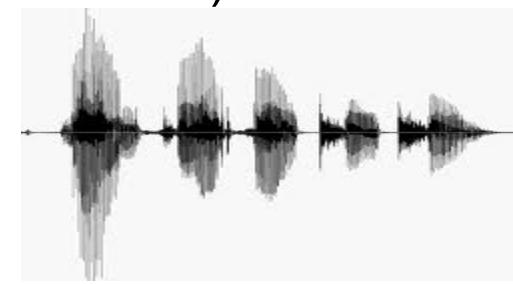
Example 3 (plots of dependences, multidimensional time-series)



Electricity price curve



Lift force vs. AoA curve



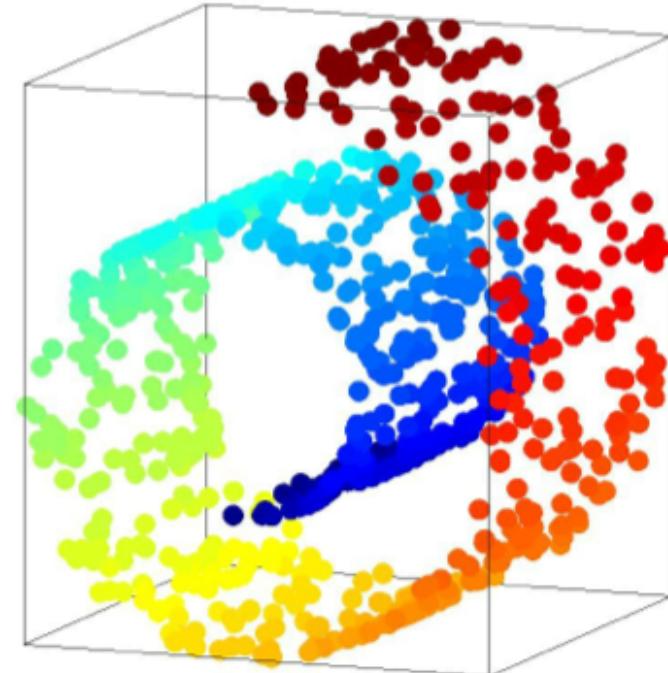
Speech recognition

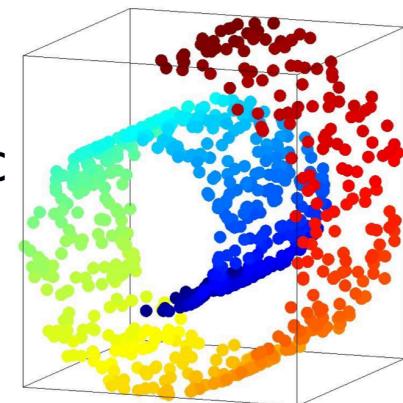
Curve $f(x)$ is described by the vector $f = (f_1, f_2, \dots, f_p)^T \in R^p$, $f_j = f(t_j)$; time-series segment (x_1, x_2, \dots, x_p) is considered as a p -dimensional vector

Example 4 (MNIST)



Example 5 Famous Toy Problem “Swiss Roll”



- High dimensionality p of $\mathbf{X}(\mathbf{O})$ is critical for efficient learning
 - We can visualize only in 2D/3D
 - Dimension Reduction = construct reduced dimension representation $\mathbf{y}(\mathbf{O}) \in \mathbb{R}^q$, $q \ll p$, of \mathbf{O} without “significant loss of information”
 - DR for
 - ✓ Visualization
 - ✓ Data compression
 - ✓ “curse of dimensionality”
 - ✓ De-noising
 - ✓ Reasonable distance metrics
- $\mathbf{X} \rightarrow \mathbf{X}'$ S.T.
 $\dim(\mathbf{X}') \ll \dim(\mathbf{X})$
- uncovers the intrinsic
dimensionality
(invertible)
- 

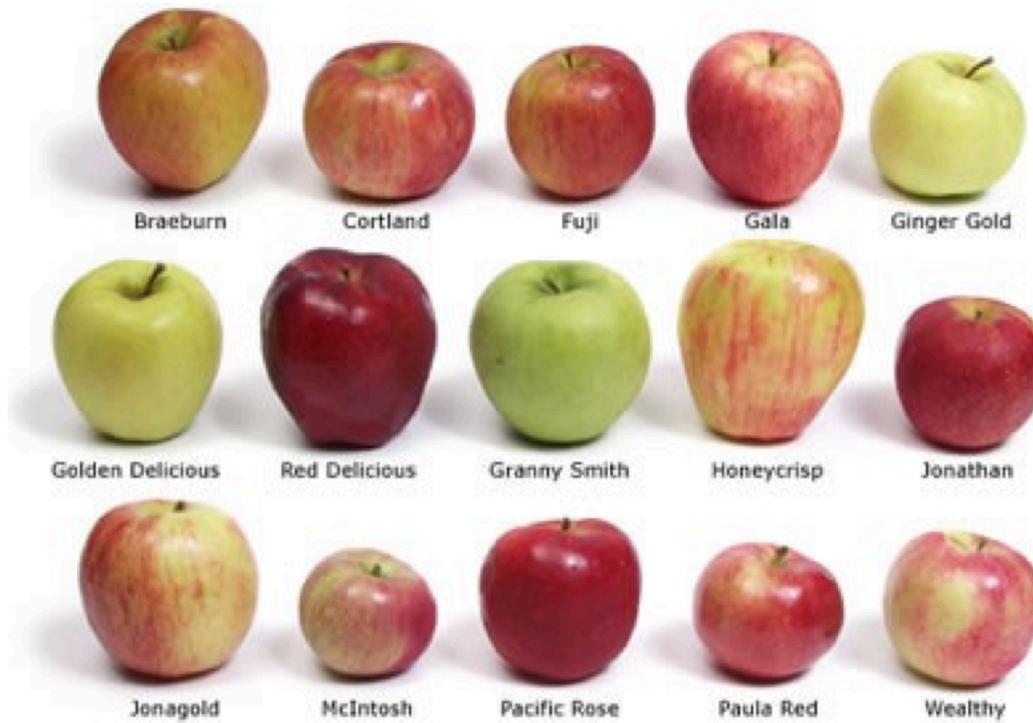
compressed description $y(O)$ is constructed for a detailed description $X(O)$ using a learning sample

$$X_n = \{X_i = X(O_i), i = 1, 2, \dots, n\} = \{X_1, X_2, \dots, X_n\}$$

of features for objects O_1, O_2, \dots, O_n ,

without using physical information about objects

appearance variation

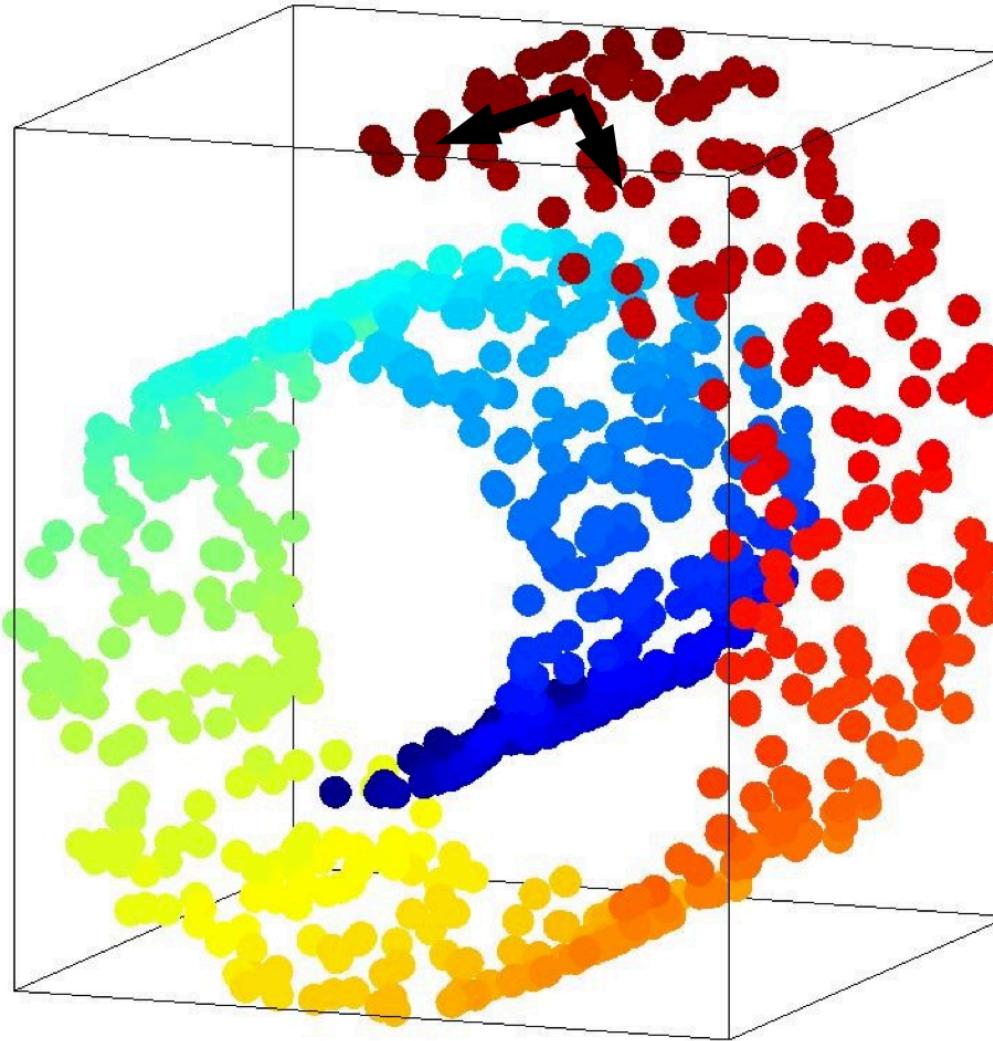


What is “Reasonable distance metrics”?

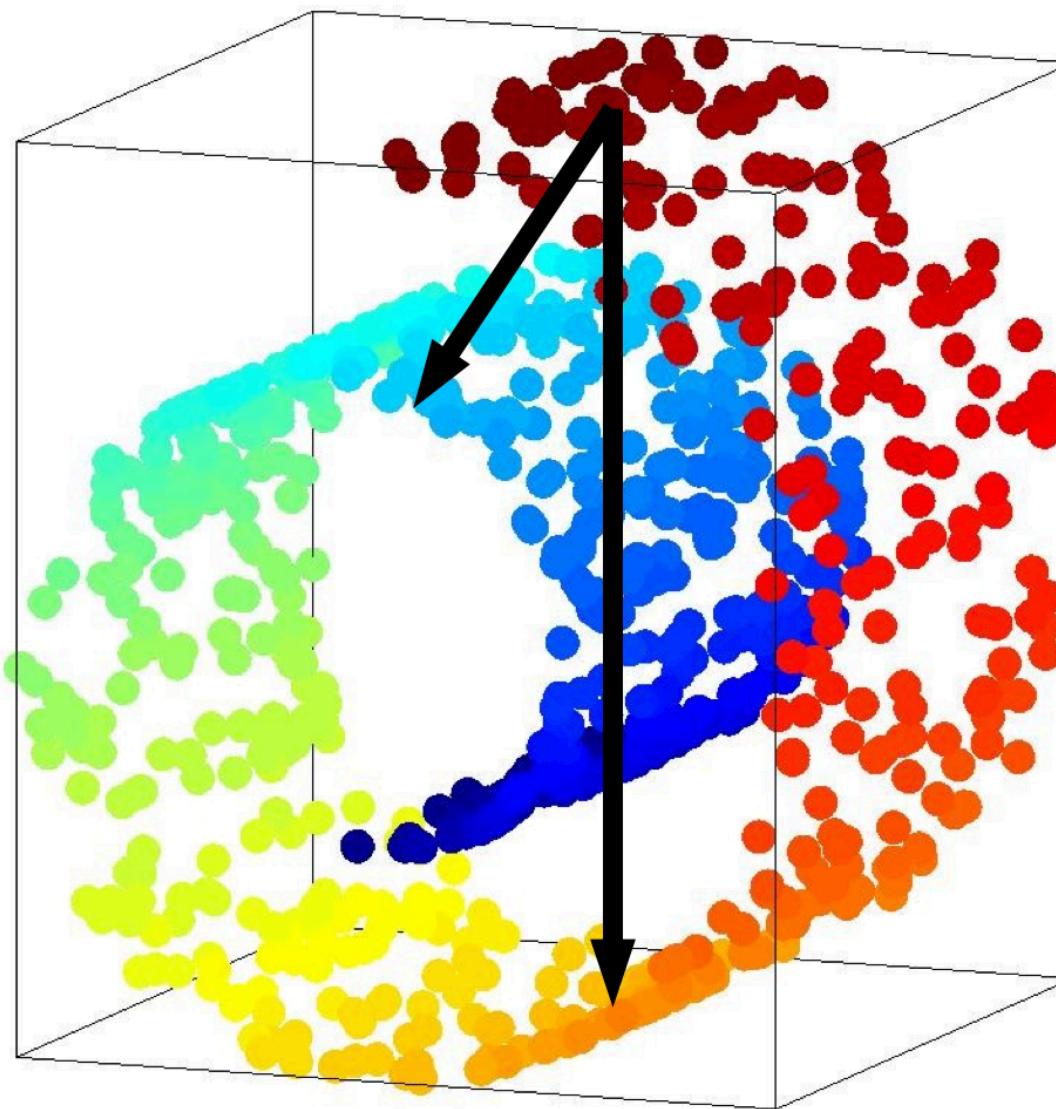
deformation



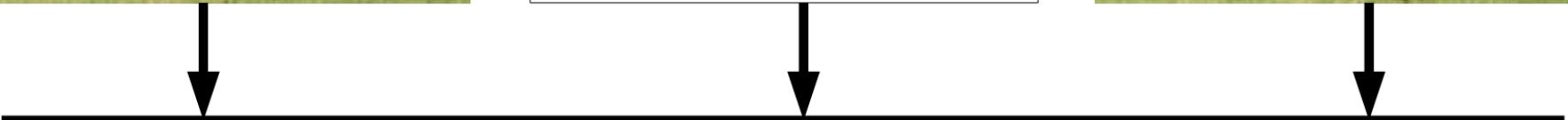
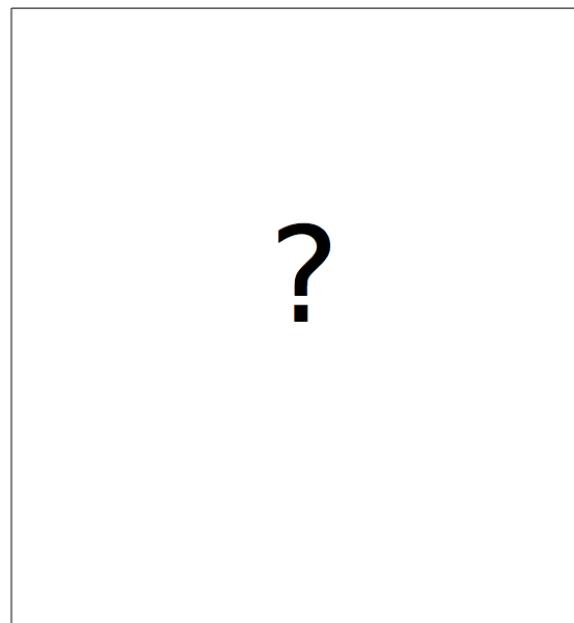
What is “Reasonable distance metrics”?



What is “Reasonable distance metrics”?



reasonable distance metrics



reasonable distance metrics



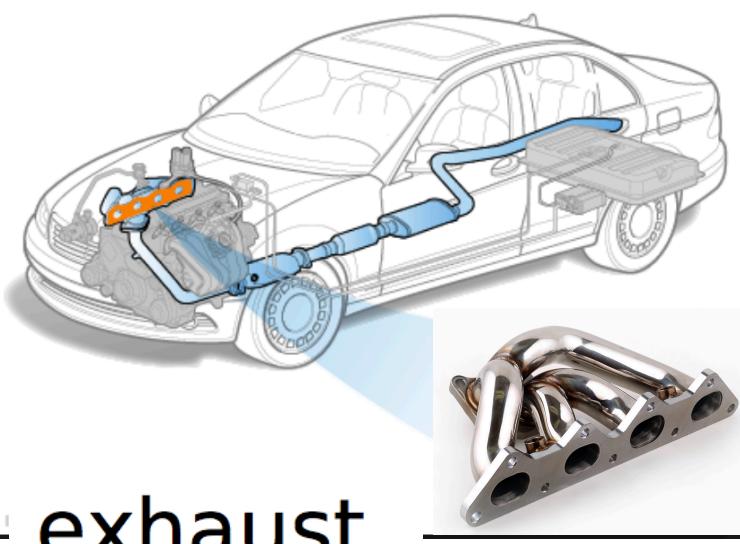
linear interpolation

reasonable distance metrics

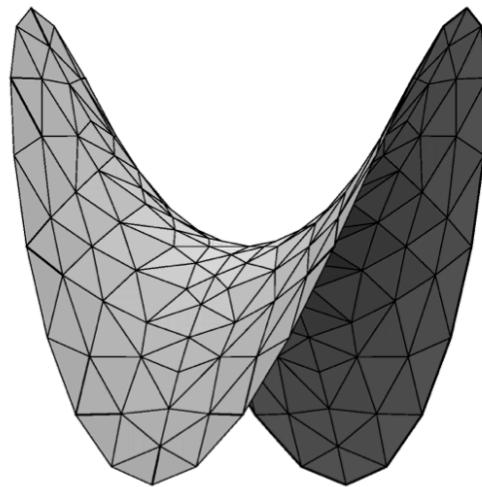


manifold interpolation

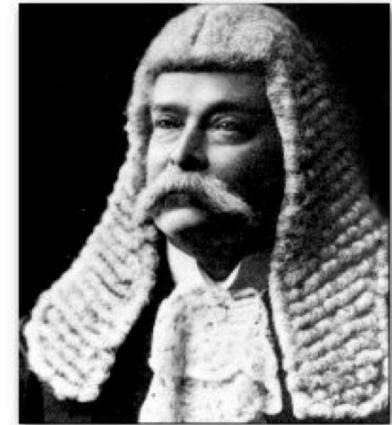
Types of manifold



- exhaust
manifold

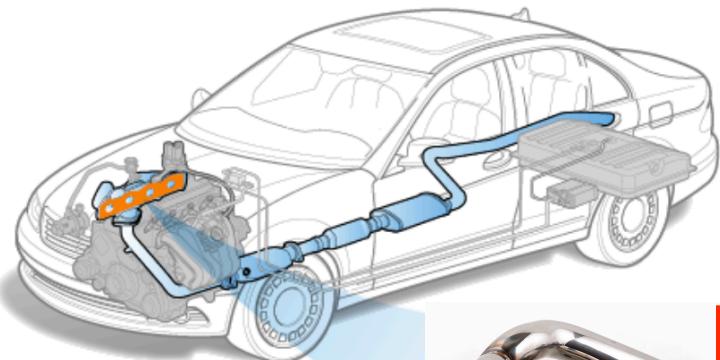


low-D surface
embedded in
high-D space

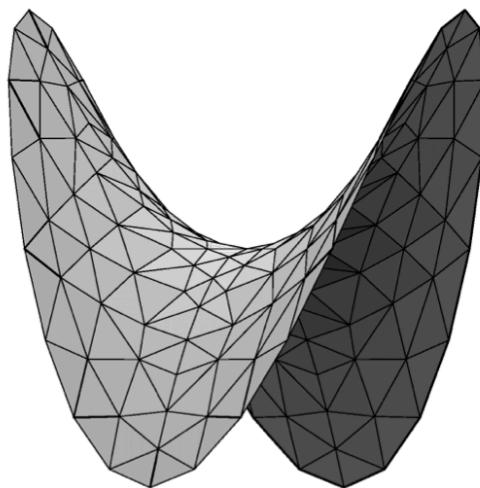


Sir Walter
Synnot Manifold
1849-1928

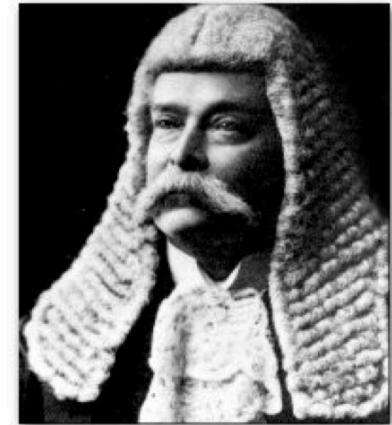
Types of manifold



- exhaust manifold



low-D surface
embedded in
high-D space



Sir Walter
Synnot Manifold
1849-1928

Dimension Reduction as an Embedding Problem

Embedding Problem. Using a sample

$$\mathbf{X}_n = \{X_1, X_2, \dots, X_n\} \subset \mathbb{R}^p$$

construct an embedding

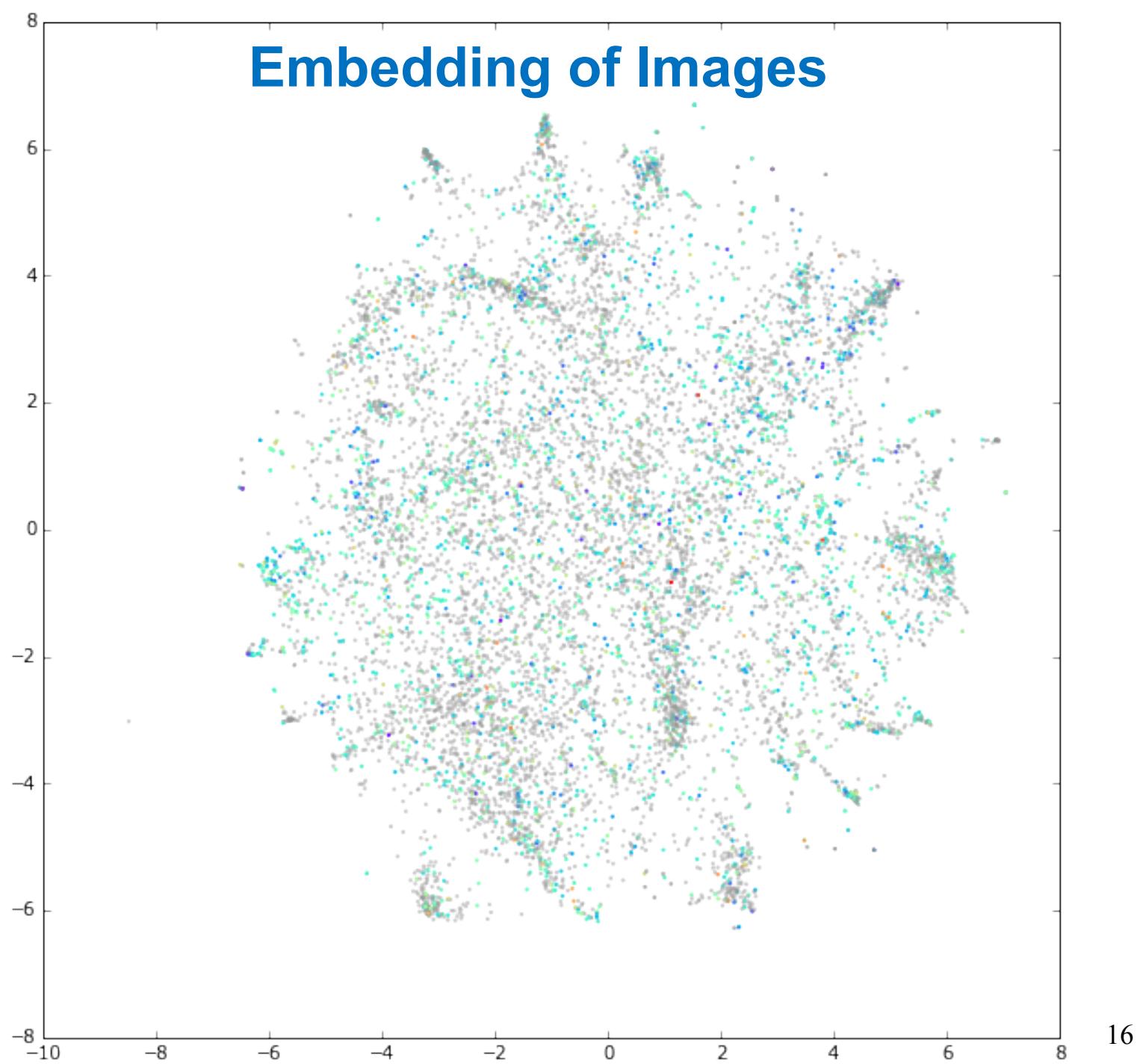
$$h: \mathbf{X}_n \rightarrow \mathbf{Y}_n = h(\mathbf{X}_n) = \{y_1, y_2, \dots, y_n\} \subset \mathbb{R}^q$$

of sample points $\mathbf{X}_n \subset \mathbb{R}^q$, $q < p$, such that the set \mathbf{Y}_n

«consistently represents» the set \mathbf{X}_n

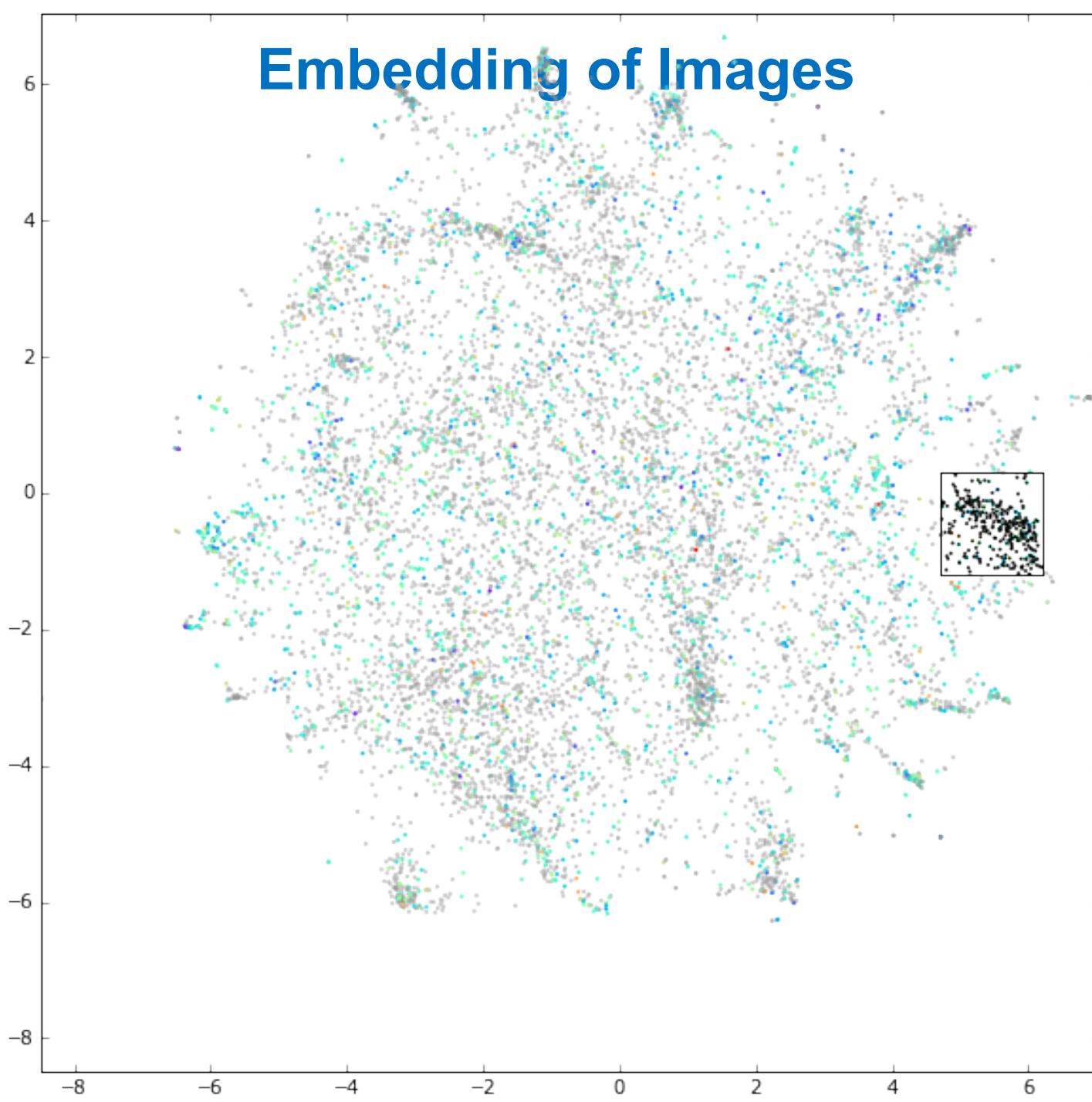
(e.g., the set \mathbf{Y}_n should preserve some geometric structure of \mathbf{X}_n , etc.)

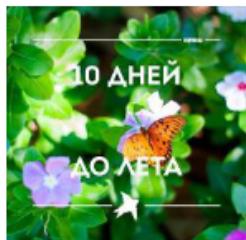
Embedding of Images



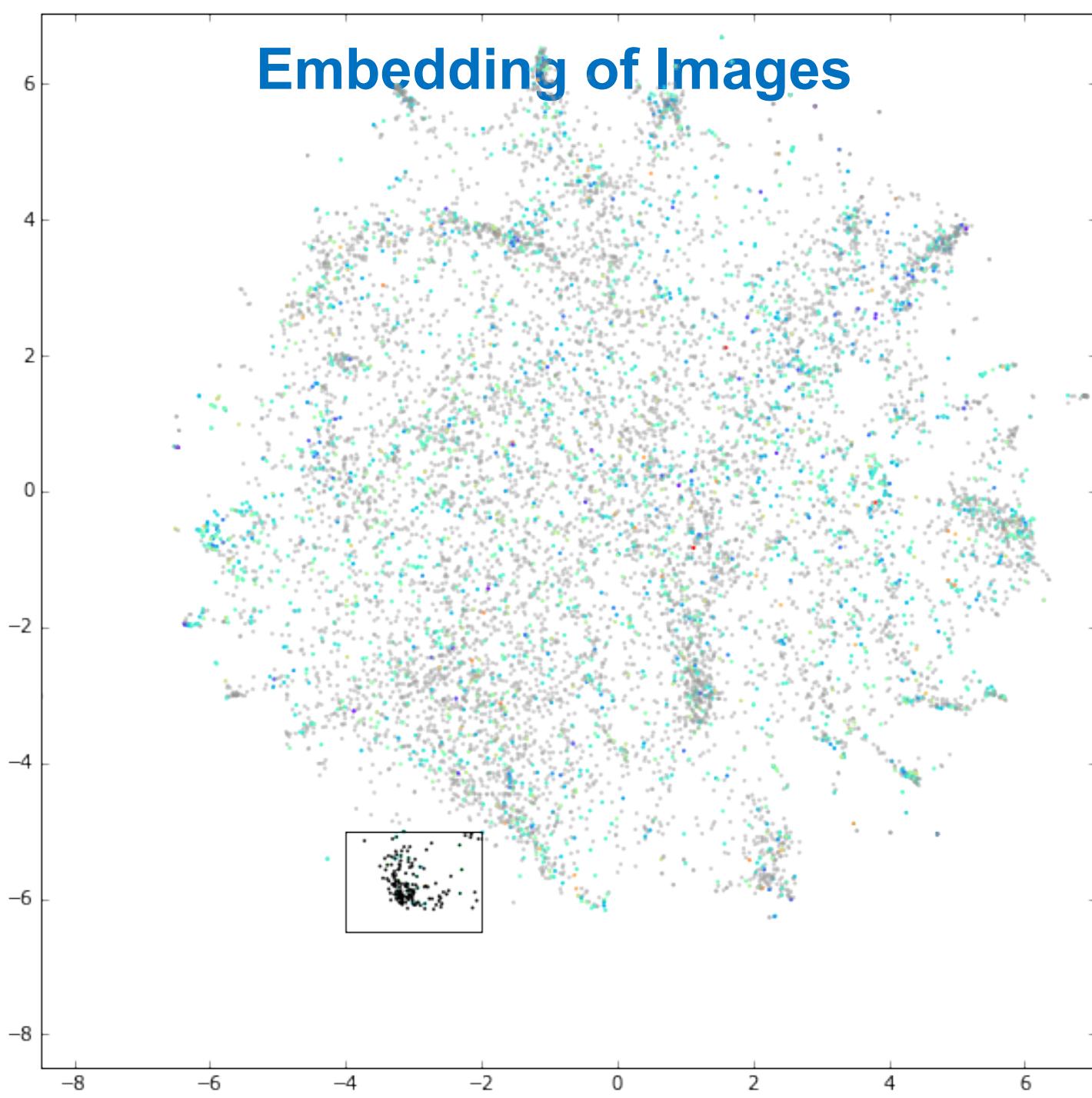
16

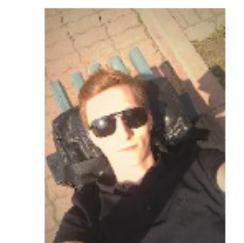
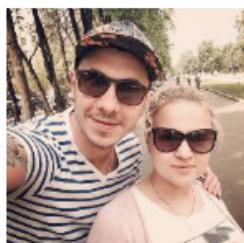
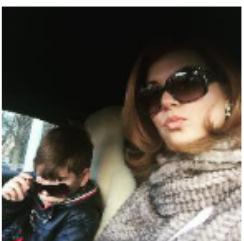
Embedding of Images



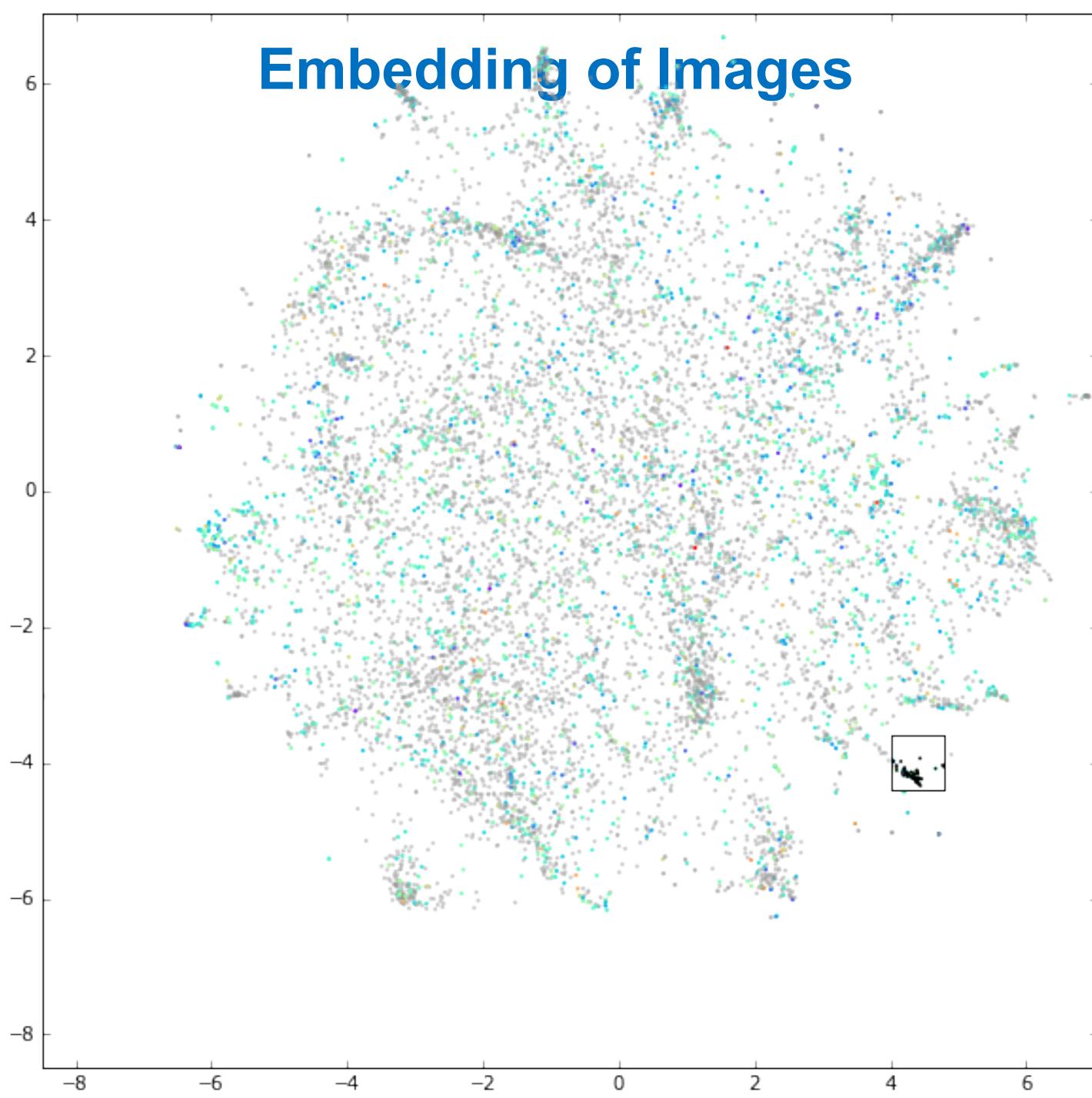


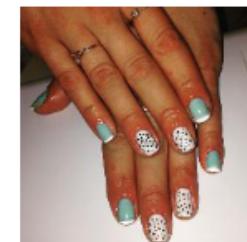
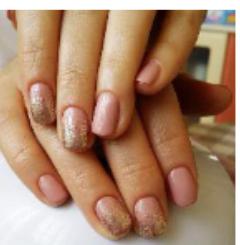
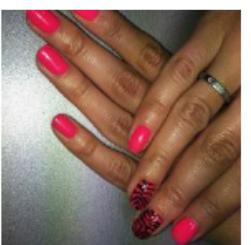
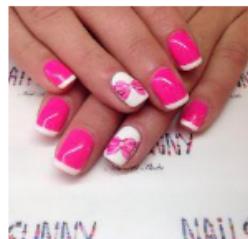
Embedding of Images





Embedding of Images





Typical scheme:

- construct some cost function $L(y_1, y_2, \dots, y_n)$. E.g. for Multi Dimensional Scaling, MDS

$$L(y_1, y_2, \dots, y_n) = \sum_{i,j} (\rho(X_i, X_j) - \|y_i - y_j\|)^2$$

- optimize cost function
- visually analyze results (Swiss Roll, Spiral, ...)

Popular methods to solve embedding problem:

- Pursuit Projection; Principal Component Analysis
- Locally Linear Embedding, LLE; Conformal Eigenmaps
- Laplacian Eigenmaps LE
- Hessian Eigenmaps, HE
- ISOfometric MAPing, ISOMAP; Landmark ISOMAP
- Kernel PCA, KPCA
- Local Tangent Space Alignment, LTSA

Extended Embedding Problem

Using a sample

$$\mathbf{X}_n = \{X_1, X_2, \dots, X_n\} \subset \mathbb{R}^p$$

construct an embedding transformation

$$h: \mathbf{X} \subset \mathbb{R}^p \rightarrow \mathbf{Y} = h(\mathbf{X}) \subset \mathbb{R}^q,$$

both for sample points X_n , and for new (out-of-sample) points

$$X_{\text{new}} \in \mathbf{X} / X_n$$

without solving embedding problem for $\mathbf{X}_n \cup \{X\}$ anew

Example – face recognition

For extended embedding problem we need a Data Model \mathbf{X} : description of \mathbf{X} , and of a mechanism, which generates points \mathbf{X}_n and new points \mathbf{X}_{new} from \mathbf{X}

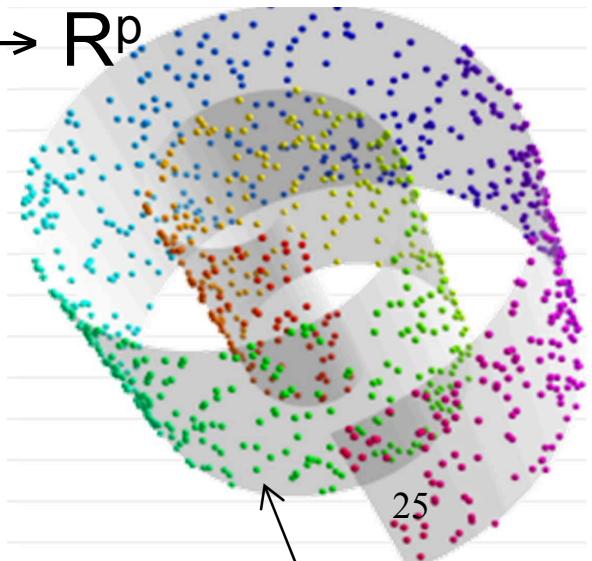
Usually we assume that real data is well approximated by some manifold, i.e.

Manifold Data Model:

$$\mathbf{X} = \{\mathbf{X} = f(\mathbf{b}) \in \mathbb{R}^p : \mathbf{b} \in \mathbf{B} \subset \mathbb{R}^q\} \subset \mathbb{R}^p$$

- open set \mathbf{B} (inner coordinate space)
- smooth bijective transformation $f: \mathbf{B} \rightarrow \mathbb{R}^p$

Manifold Learning Problem!!!



Full Dimension Reduction problem

Using a sample

$$\mathbf{X}_n = \{X_1, X_2, \dots, X_n\} \subset \mathbb{R}^p$$

construct embedding transformation:

$$h: \mathbf{X} \subset \mathbb{R}^p \rightarrow \mathbf{Y} = h(\mathbf{X}) \subset \mathbb{R}^q$$

and **Reconstruction transformation**

$$g: \mathbf{Y} \subset \mathbb{R}^q \rightarrow \mathbf{X} \subset \mathbb{R}^p$$

such that

$$g(h(X)) \approx X \text{ for all } X \in \mathbf{X}$$

Example – dimension reduction of airfoils

Dimension Reduction

Full description \mathbf{X}
of dimension p

DR procedure

Compressed
description $\mathbf{h}(\mathbf{X})$
of dimension q

Compression procedure (\mathbf{h})

Reconstruction procedure (\mathbf{g})

\mathbf{X}

$\mathbf{h}(\mathbf{X})$

$\mathbf{X}^* = \mathbf{g}(\mathbf{h}(\mathbf{X}))$

Requirements:

minimal dimension $\mathbf{q} = \text{Dim } \mathbf{h}(\mathbf{X})$ of compressed description $\mathbf{h}(\mathbf{X})$,
providing required proximity between $\mathbf{X}^* = \mathbf{g}(\mathbf{h}(\mathbf{X}))$ and \mathbf{X}

or:

maximal proximity between \mathbf{X} and $\mathbf{X}^* = \mathbf{g}(\mathbf{h}(\mathbf{X}))$ for a fixed dimension
 $\mathbf{q} = \text{Dim } \mathbf{h}(\mathbf{X})$ of compressed description

Why do we need reconstruction?

Optimize $v(X)$ w.r.t. p-dimensional vector $X \in X \subset R^p$. If $\theta = (h, g)$ is such that $X \approx g(h(X))$, and

$$v(X) \approx v(g(h(X))),$$

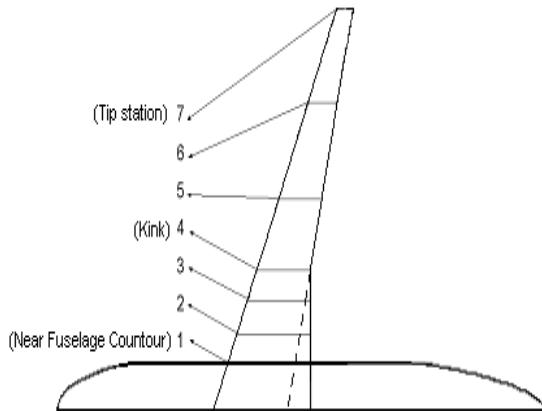
then instead we can optimize the functional

$$v_\theta(y) = v(g(y))$$

w.r.t. q-dimensional vector y

For an optimal y^* we have to recover optimal $X^* = g(y^*)$

Example – aircraft wing shape optimization



- airfoil description dimension $p = 60$
- compressed description dimension $q = 6$
- dimension for the initial problem
$$X = (X_1, X_2, \dots, X_7) \in R^{420}$$
- dimension of the reduced problem
$$y = (y_1, y_2, \dots, y_7) \in R^{42}$$
- For optimal y we have to recover a real wing description

Principal Component Analysis, PCA

(K. Pearson, 1899)

Problem: find in \mathbb{R}^p an affine subspace

$$L(q) = \left\{ x \in \mathbb{R}^p : x = x_0 + \sum_{j=1}^q y_j \times e_j, y_1, y_2, \dots, y_q \in \mathbb{R}^1 \right\}$$

of dimension $q < p$, **which best approximates the set of points**

$$\mathbf{X}_n = \{X_i, i = 1, 2, \dots, n\} \subset \mathbb{R}^p.$$

in PCA: “the best” = minimize w.r.t. $\mathbf{x}_0, \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_q\} \subset \mathbb{R}^p$

$$\frac{1}{n} \sum_{j=1}^n \|X_j - P_{L(q)} X_j\|^2$$

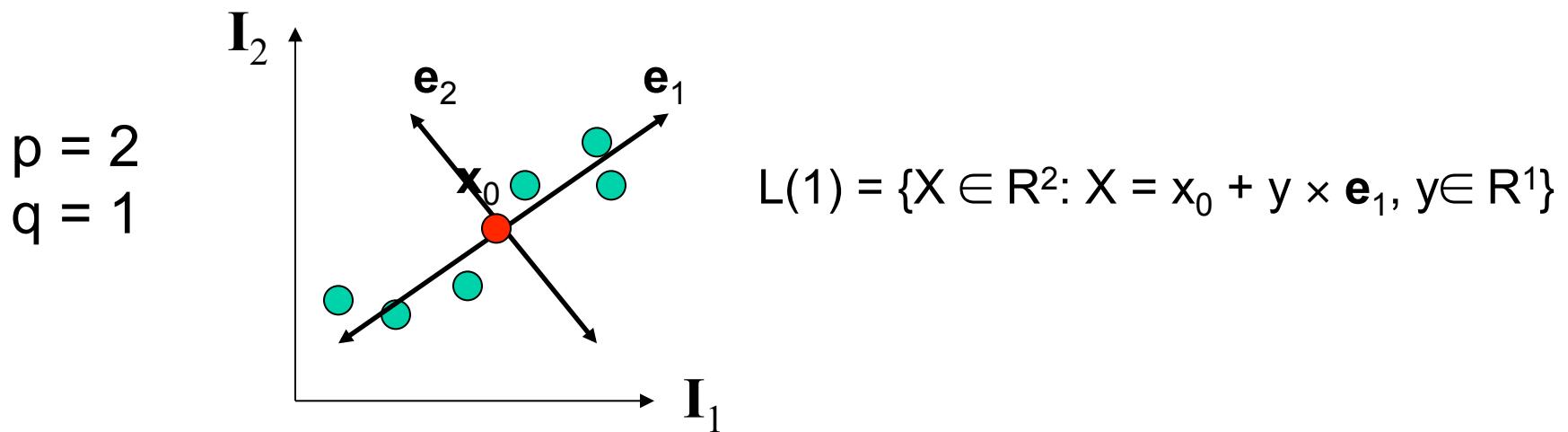
$$\text{Pr}_{L(q)}(X) = x_0 + \sum_{j=1}^q y_j(X) \times e_j, \quad y_j(X) = (X - x_0, e_j)$$

x_{mean} – empirical mean of $\{X_i, i = 1, 2, \dots, n\}$,

$\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p\}$ – eigenvalues of an empirical covariance ($p \times p$)-matrix

$$\Sigma = \frac{1}{n} \sum_{j=1}^n (X_j - x_{\text{mean}}) \times (X_j - x_{\text{mean}})^T$$

providing orthonormal basis in \mathbb{R}^p



Solution: $x_0 = x_{\text{mean}}$, $L(q) = x_0 \oplus \text{Span}(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_q)$,

where orthonormal eigenvectors $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_q\}$ of Σ

correspond to q highest eigenvalues of this matrix

PCA solves E-Problem, EE-Problem and FULL DR problem:

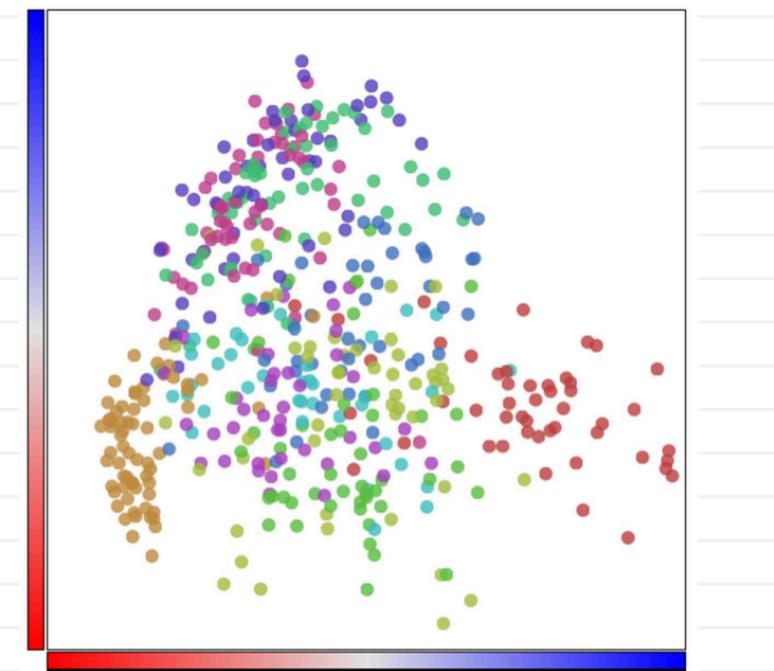
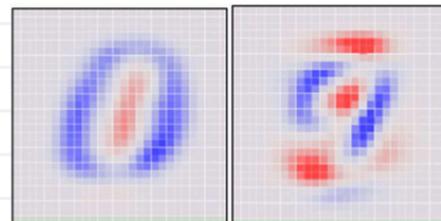
embedding procedure:

$$X \in \mathbb{R}^p \rightarrow h(X) = (y_1(X), y_2(X), \dots, y_q(X))^T \in \mathbb{R}^q$$

reconstruction procedure:

$$y = (y_1, y_2, \dots, y_q)^T \in \mathbb{R}^q \rightarrow g(y) = x_0 + y_1 \times e_1 + y_2 \times e_2 + \dots + y_q \times e_q$$

PCA on MNIST



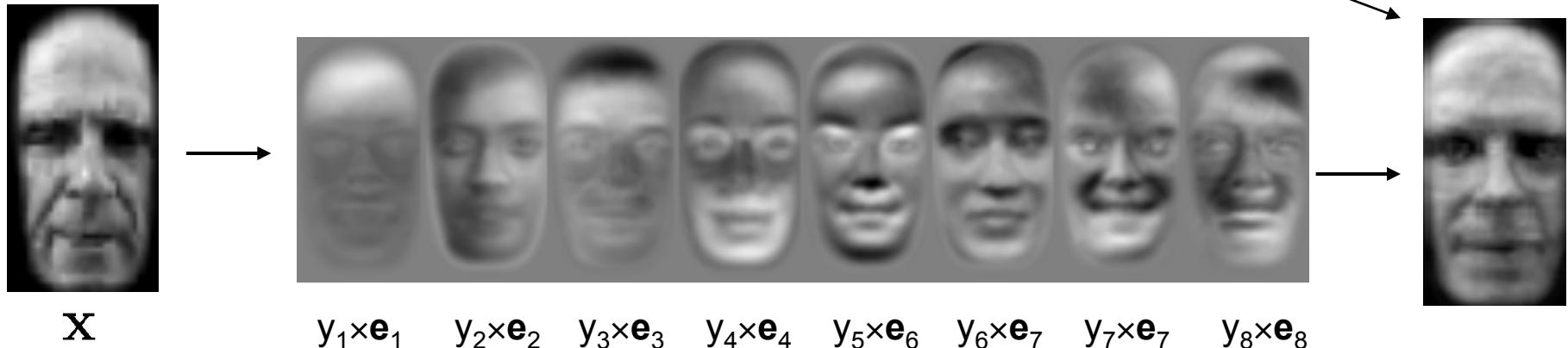
PCA for Face Recognition

- Eigenvectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p \in \mathbb{R}^p$ in a space of faces: $p \sim 10^6$

$$X \rightarrow ((\underbrace{(X - x_0, \mathbf{e}_1)}_{y_1}), (\underbrace{(X - x_0, \mathbf{e}_2)}_{y_2}), \dots, (\underbrace{(X - x_0, \mathbf{e}_p)}_{y_p}))^\top \in \mathbb{R}^p$$

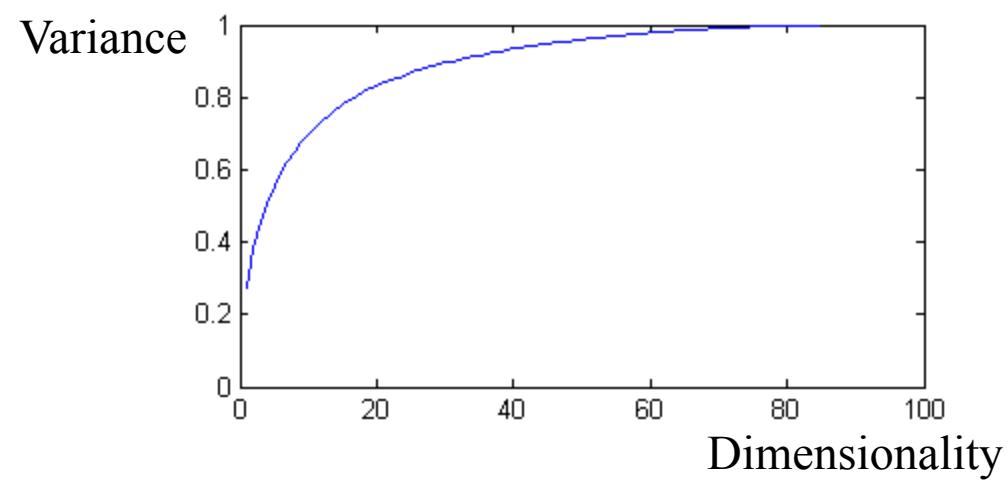
- We select first q , $q \sim 10^2$ eigenvectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_q \in \mathbb{R}^q$

$$X^* \approx x_0 + y_1 \times \mathbf{e}_1 + y_2 \times \mathbf{e}_2 + \dots + y_q \times \mathbf{e}_q$$





Left to right: original, reconstruction from 84, 40, 20, 3, 2, and 1 dimensions.



Pursuit Projection, PP

Similar to PCA, constructs «**the best**» affine hyperplane of smaller dimension

but for some other performance criteria

E.g. to better highlight clusters in data

One of the first papers is by J.H. Friedman and J. W. Tukey

“A Projection Pursuit Algorithm for Exploratory Data Analysis” (1974)

Multi Dimensional Scaling, MDS

Minimize quadratic form

$$\sum_{i,j} (\rho(O_i, O_j) - \|y_i - y_j\|)^2$$

w.r.t. $y_1, y_2, \dots, y_n \in R^q$, ρ is a some metrics in the object feature space

$y_1, y_2, \dots, y_n \in R^q$ are defined **except for shift and rotation**, thus we use normalization, e.g.

$$Y^T \times Y = I_q \quad \text{and} \quad Y^T \times 1 = 0$$

where $Y^T = (y_1 : y_2 : \dots : y_n)$ – $(q \times n)$ -matrix, $1 \in R^n$ – vectors of ones

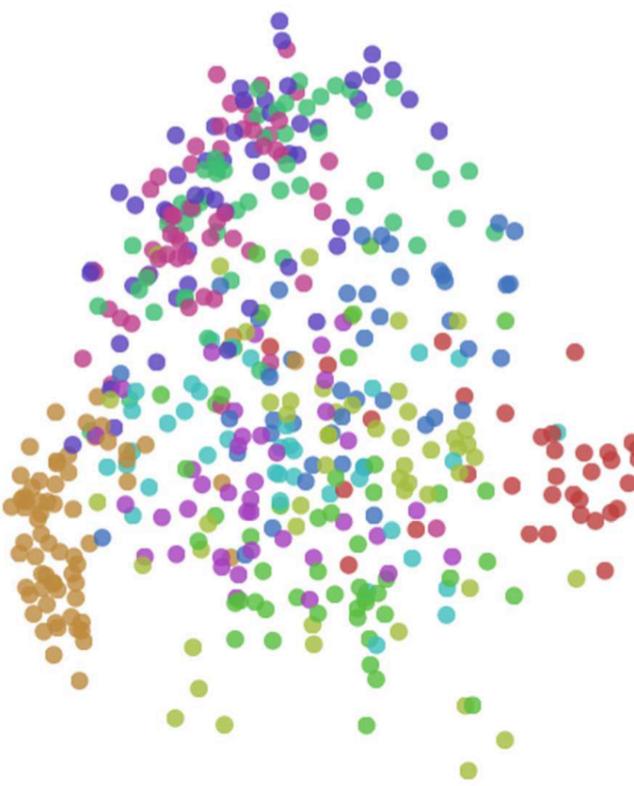
If $\rho(O_i, O_j) = \|X(O_i) - X(O_j)\|$, then MDS = PCA

Sammon mapping

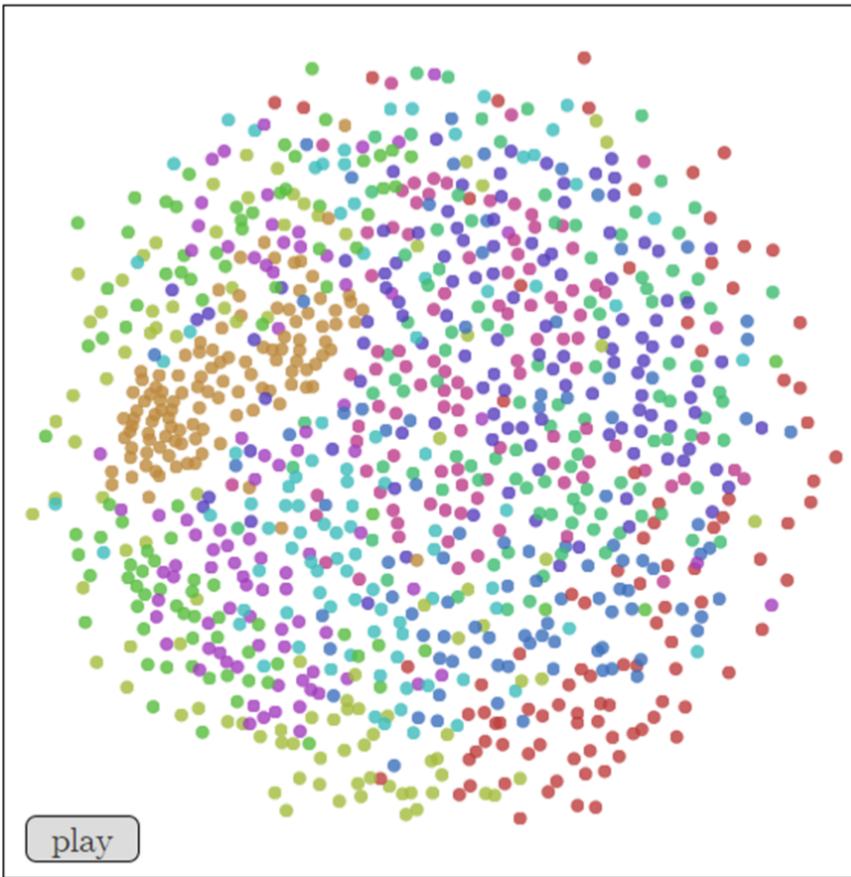
Common useful idea: local distances are often more important to preserve

$$\sum_{i,j} \frac{(\rho(O_i, O_j) - \|y_i - y_j\|)^2}{(\rho(O_i, O_j))^2}$$

PCA vs. MDS

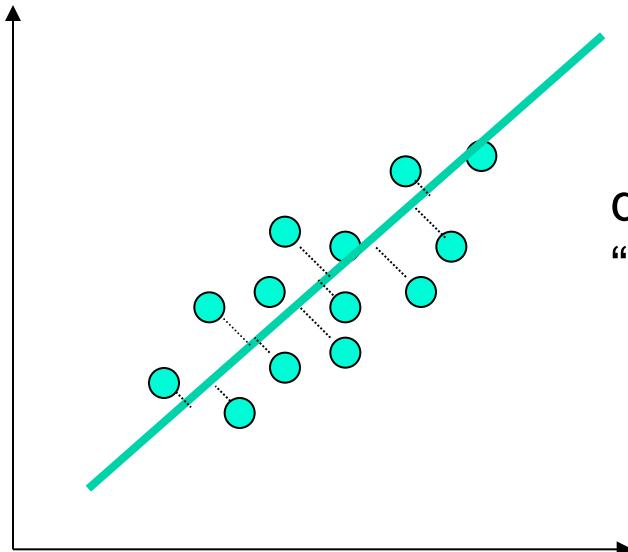


PCA

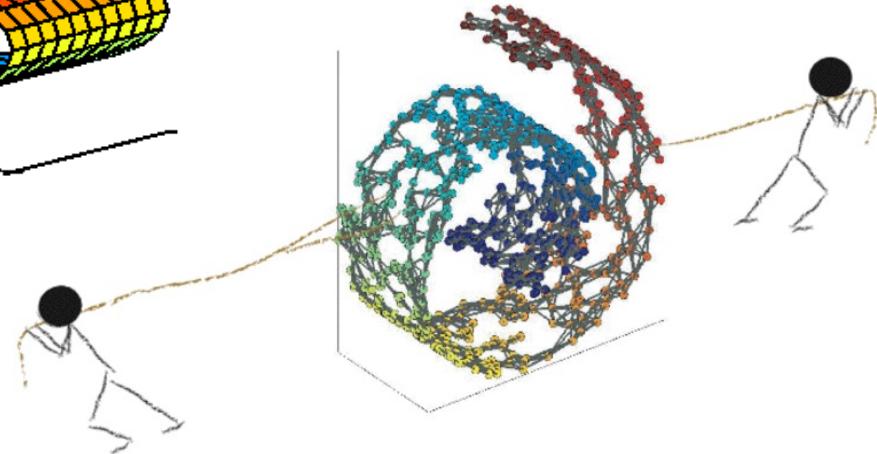
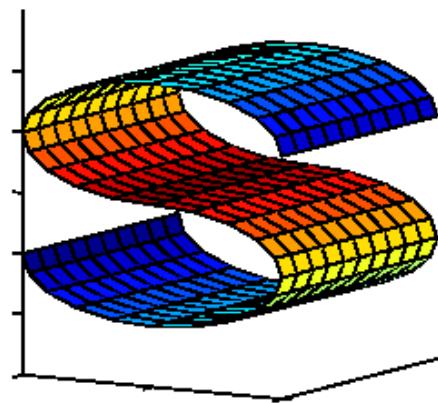
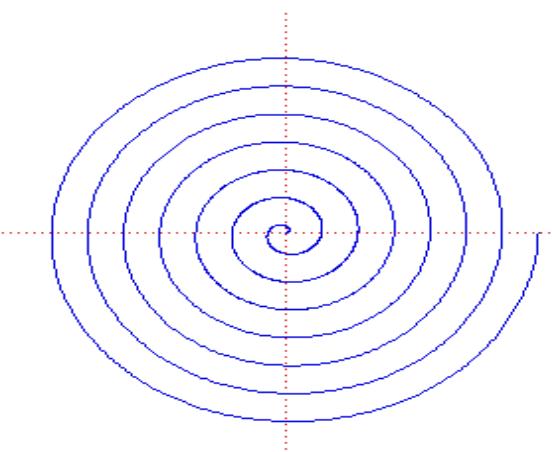
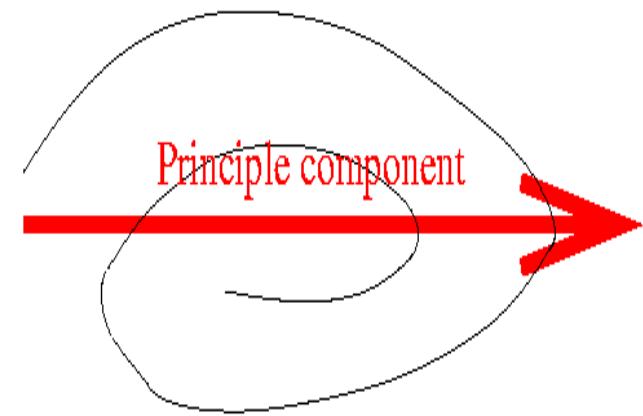


Sammon

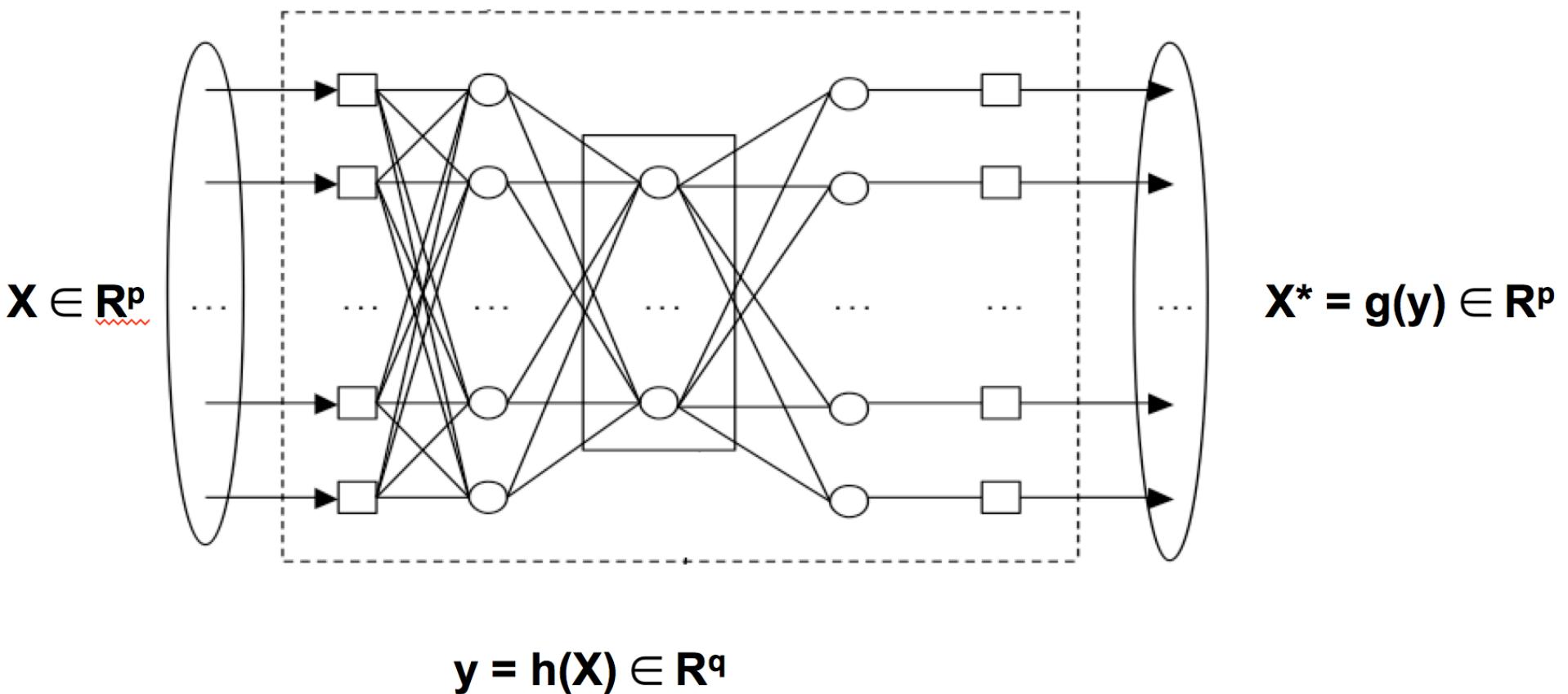
MDS, PCA, PP, MDS (for «Euclidean proximity») – linear methods,

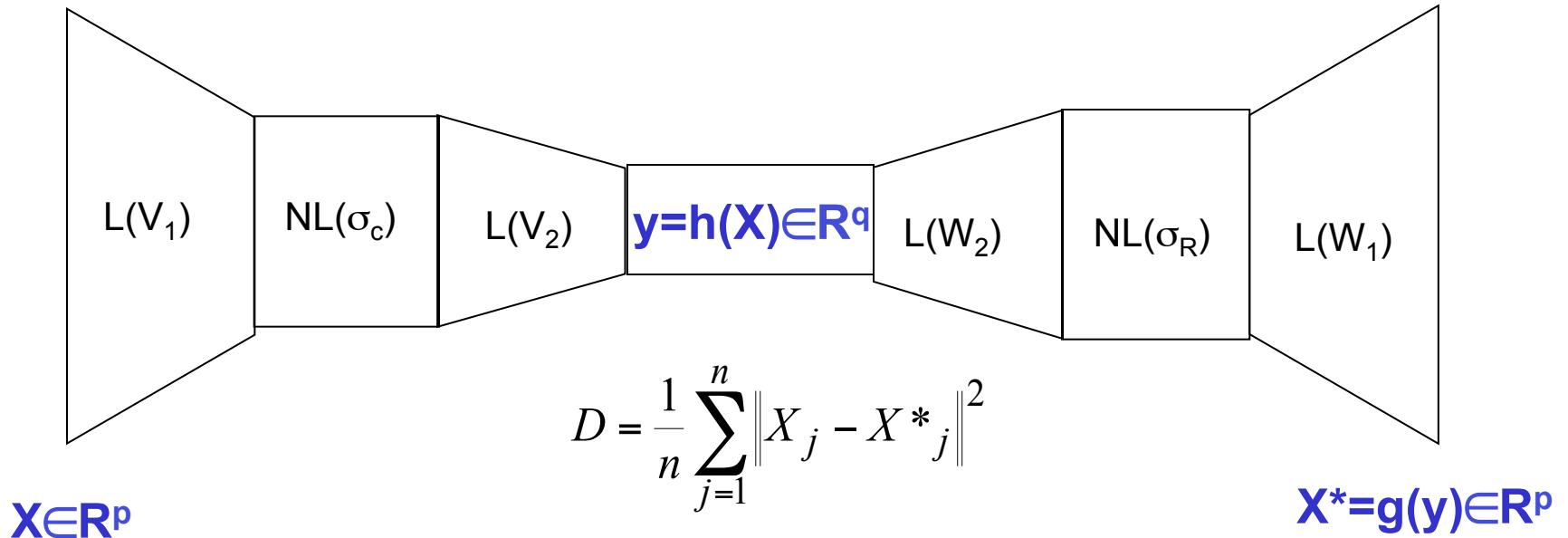


do not work for
“nonlinear data”:



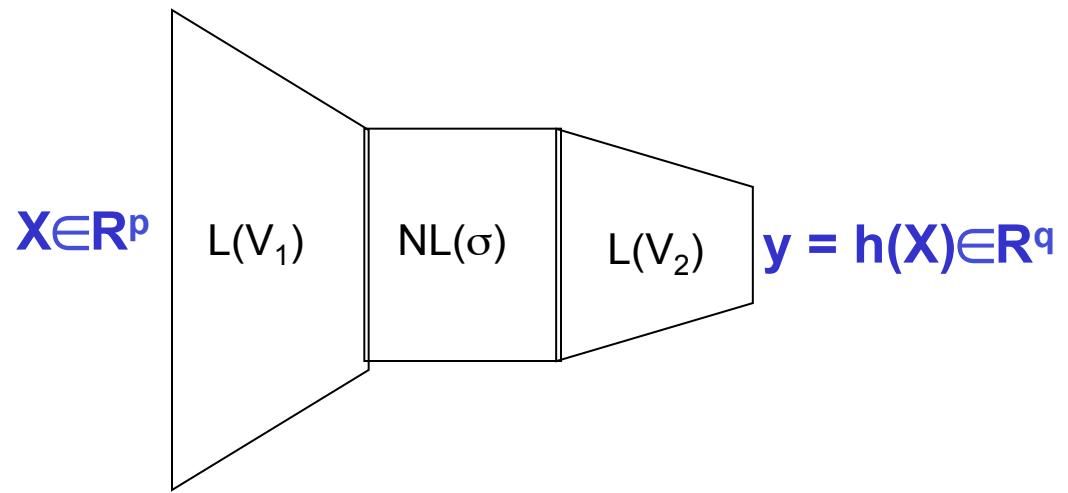
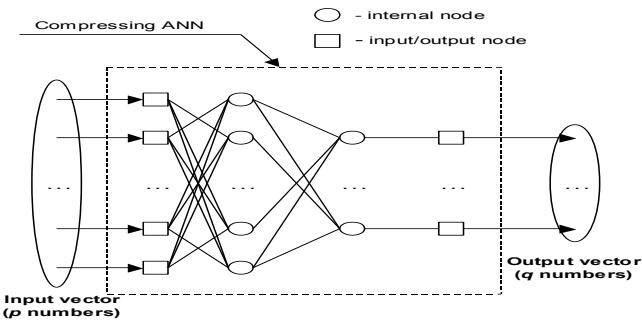
Replicative Neural Networks



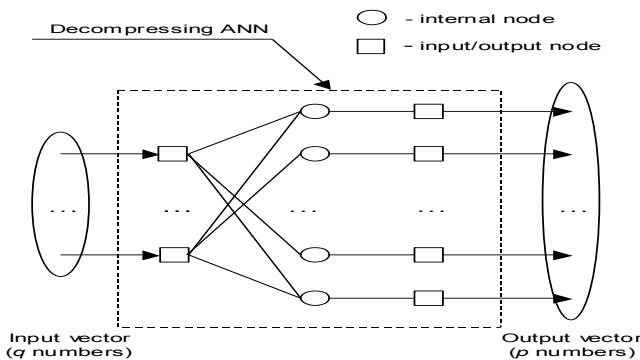


- 4 linear transformation with matrices V_1 , V_2 , W_1 and W_2
- 2 **fixed nonlinear transformations**
 $\sigma(x)$ – **sigmoid functions**, e.g., $\sigma(x) = (1 + e^{-x})^{-1}$
- Back-Propagation to optimize w.r.t. to parameters V_1 , V_2 , W_1 and W_2 the reconstruction error D

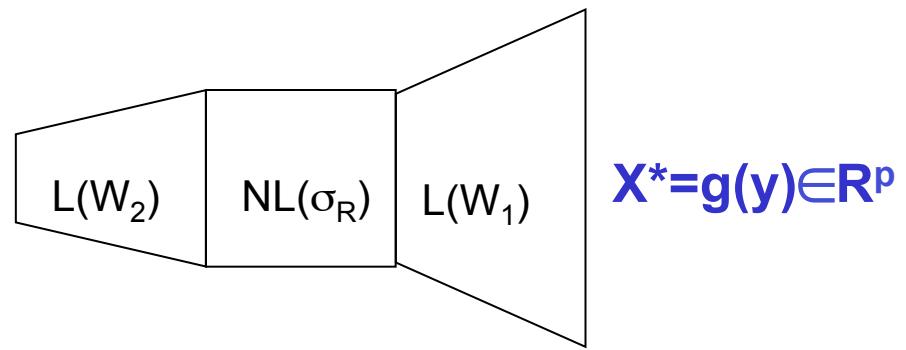
Compression transformation



Reconstruction transformation



$$y \in \mathbb{R}^q$$

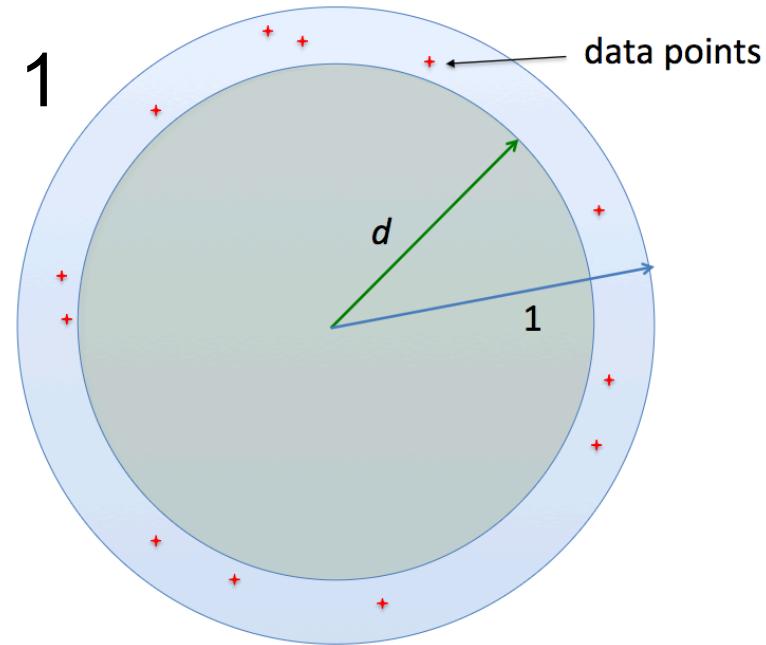
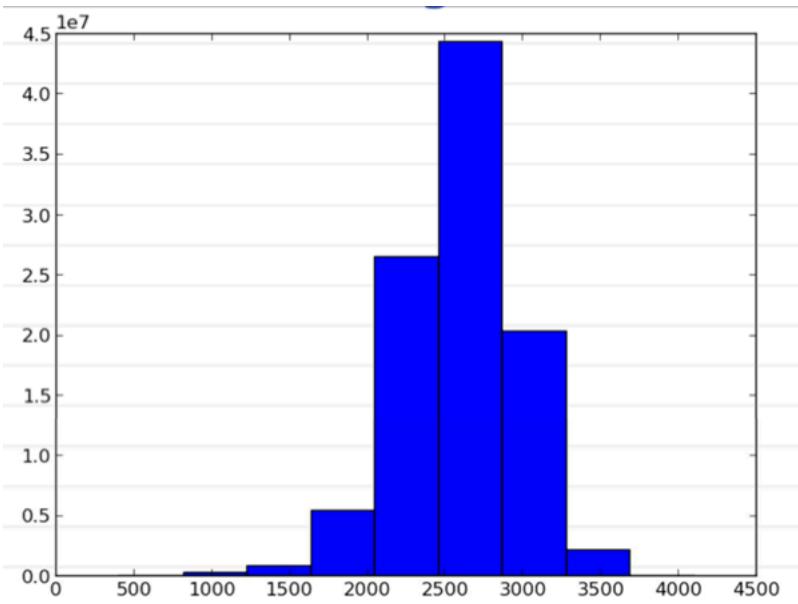


Curse of dimensionality

$$V(r) \sim C(d) * r^d$$

All volume of a ball is close to its cover

$$\frac{V(r) - V(r(1-\varepsilon))}{V(r)} = \frac{1 - (1-\varepsilon)^d}{1} \rightarrow 1$$



Histogram of pairwise Euclidean distances for MNIST

=> Consider more clever distance!!!!

Nonlinear (local) DR methods

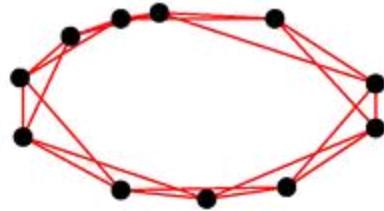
- 5) Locally Linear Embedding, LLE**
- 6) Laplacian Eigenmaps LE**
- 7) Hessian Eigenmaps, HE**
- 8) ISOmetric MAPing, ISOMAP**
- 9) Kernel PCA, KPCA - Spectral Embedding Algorithm, SEA**
- 10) Riemannian Manifold Learning, RML**
- 11) Local Tangent Space Alignment, LTSA, ...**

GRAPH: standard step for many “local” DR methods

Step 1. Construct neighborhoods (ε -Neighborhoods, k Nearest Neighbors)

$$U(X) = U(X|\rho, \varepsilon) = \{X_i \in X_n : \rho(X, X_i) \leq \varepsilon\}$$

\Rightarrow Graph $\Gamma(X_n) = (X_n, V)$: $(X_i, X_j) \in V \Leftrightarrow X_j \in U(X_i) \text{ and } X_i \in U(X_j)$



Step 2. Construct weights $w_{ij} = w(X_i, X_j) = w(v)$, $v = (X_i, X_j) \in V$

$w_{0,ij} = 1$ for all $(X_i, X_j) \in V$, $w_{0,ij} = 0$ in other cases

ISOMAP

Usual MDS

For n data points, and a distance matrix D ,

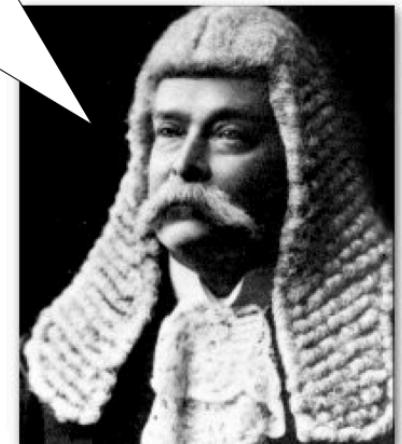
$$D_{ij} = \begin{bmatrix} & i \\ & \downarrow \\ & \rightarrow \\ & j \end{bmatrix}$$

...we can construct a m -dimensional space to preserve inter-point distances by using the top eigenvectors of D scaled by their eigenvalues.

$$y_i = [\sqrt{\lambda_1} v_1^i, \sqrt{\lambda_2} v_2^i, \dots, \sqrt{\lambda_m} v_m^i]$$

ISOMAP

Infer a distance matrix using distances along the manifold.

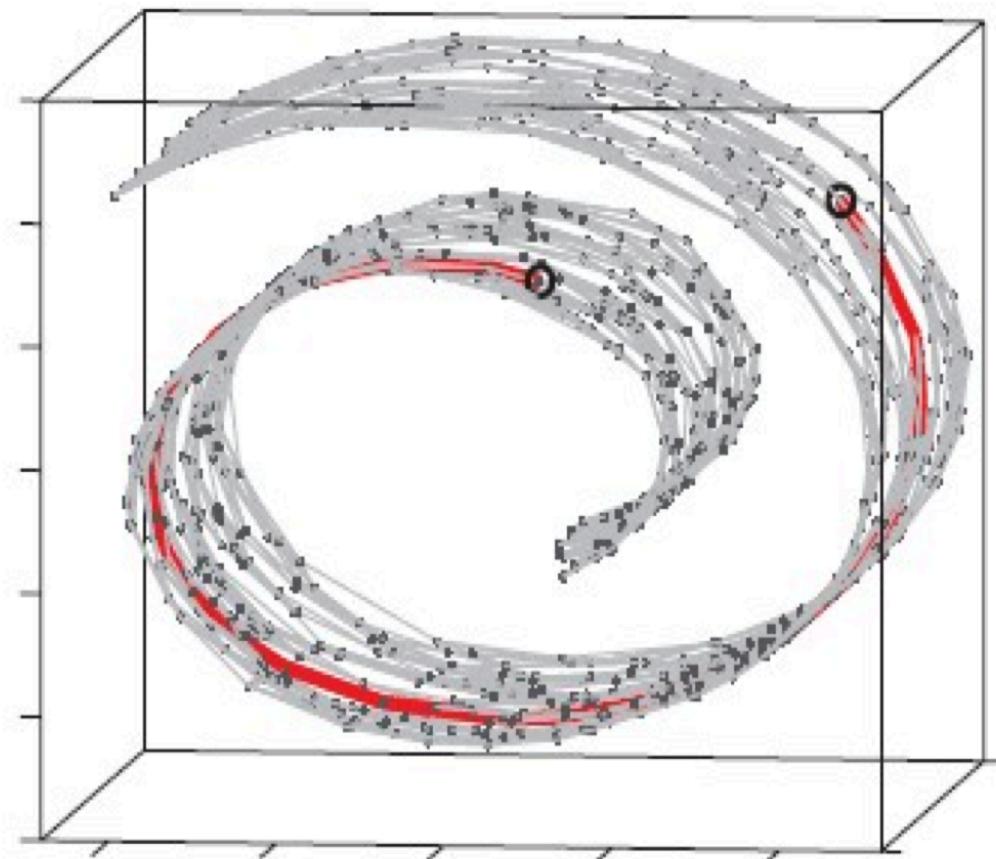


ISOMAP

1. Build a sparse graph with K-nearest neighbors

$$D_g = \begin{bmatrix} & \\ & \text{blue oval} \\ & \end{bmatrix}$$

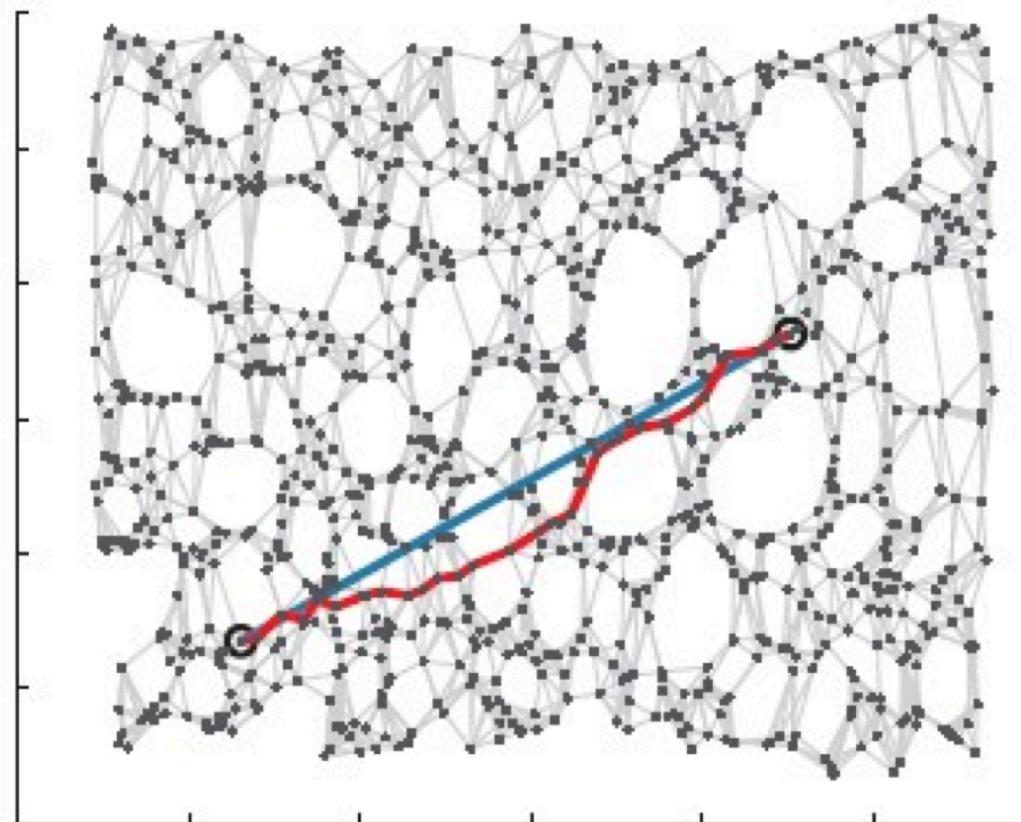
(distance matrix is sparse)



ISOMAP

2. Infer other interpoint distances by finding shortest paths on the graph (Dijkstra's algorithm).

$$D_g = \begin{bmatrix} & \\ & \end{bmatrix}$$



ISOMAP

Usual MDS

3. Build a low-D embedded space to best preserve the complete distance matrix.

Error function:

$$E = \|\tau(D_G) - \tau(D_Y)\|_{L^2}$$

inner product
distances in new
coordinate
system

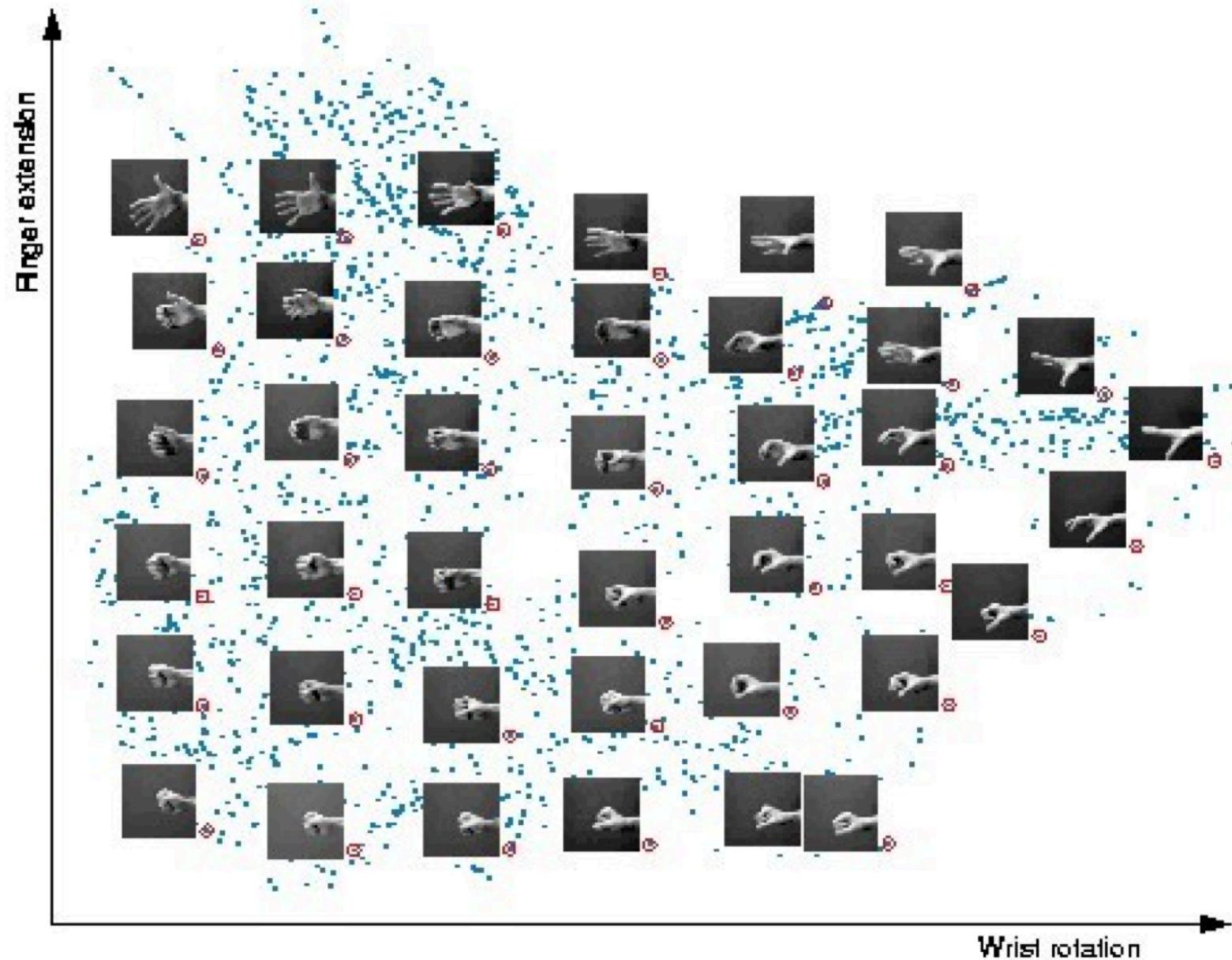
inner product
distances in graph

L2 norm

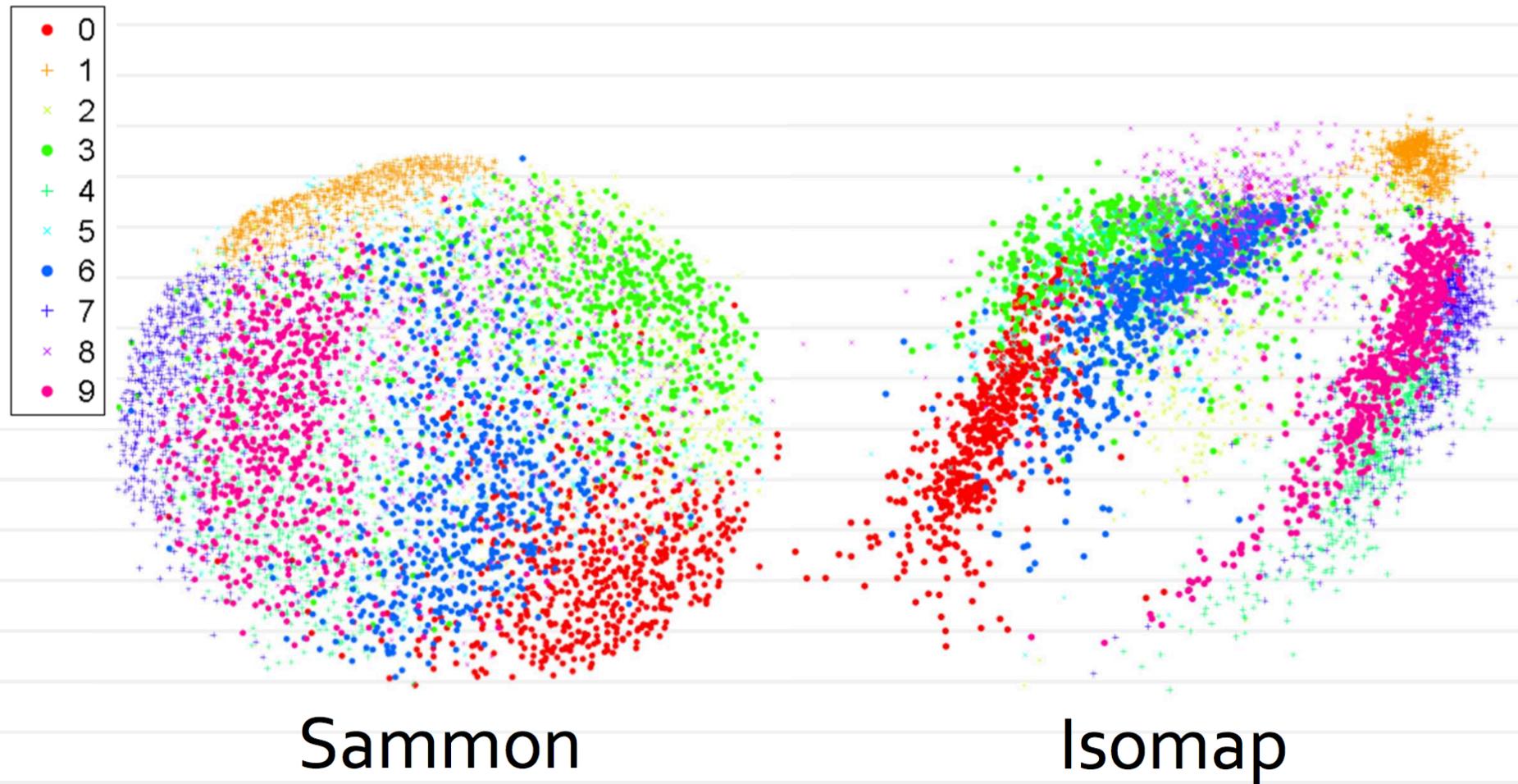
```
graph TD; E["E = ||tau(D_G) - tau(D_Y)||_{L^2}"] --> D_G["inner product distances in graph"]; E --> D_Y["inner product distances in new coordinate system"]; E --> L2_norm["L2 norm"]
```

Solution – set points Y to top eigenvectors of D_g

Isomap results: hands



MNIST: Sammon vs. ISOMAP



Sammon

Isomap

Isomap: pro and con

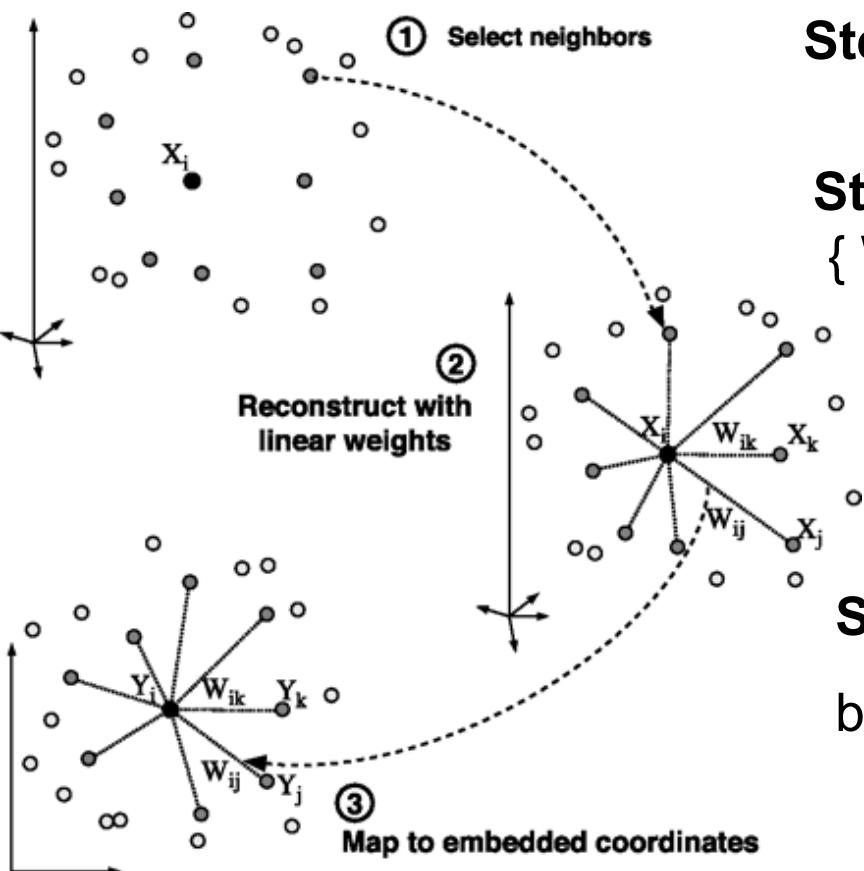
- preserves global structure
- few free parameters
- sensitive to noise, noise edges
- computationally expensive (dense matrix eigen-reduction)

Locally Linear Embedding

Find a mapping to preserve
local linear relationships
between neighbors



Locally Linear Embedding, LLE (L.K. Saul, et al., 2000)



Step 1. K nearest neighbors

Step 2. Get “Baricentric” coordinates $\{W_{ji}\}$ by minimizing

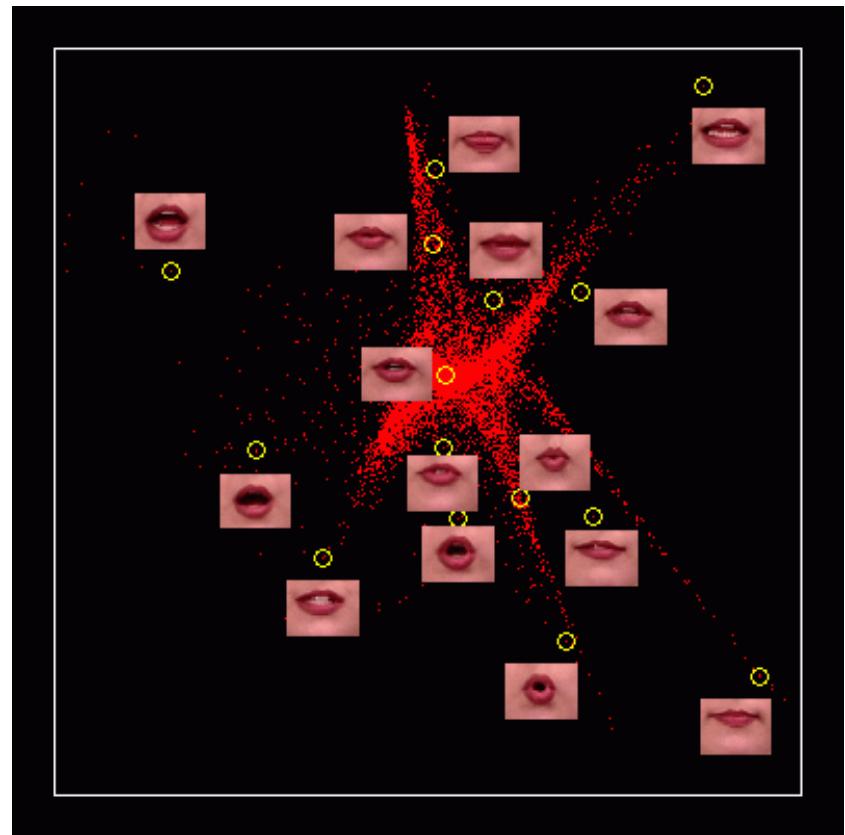
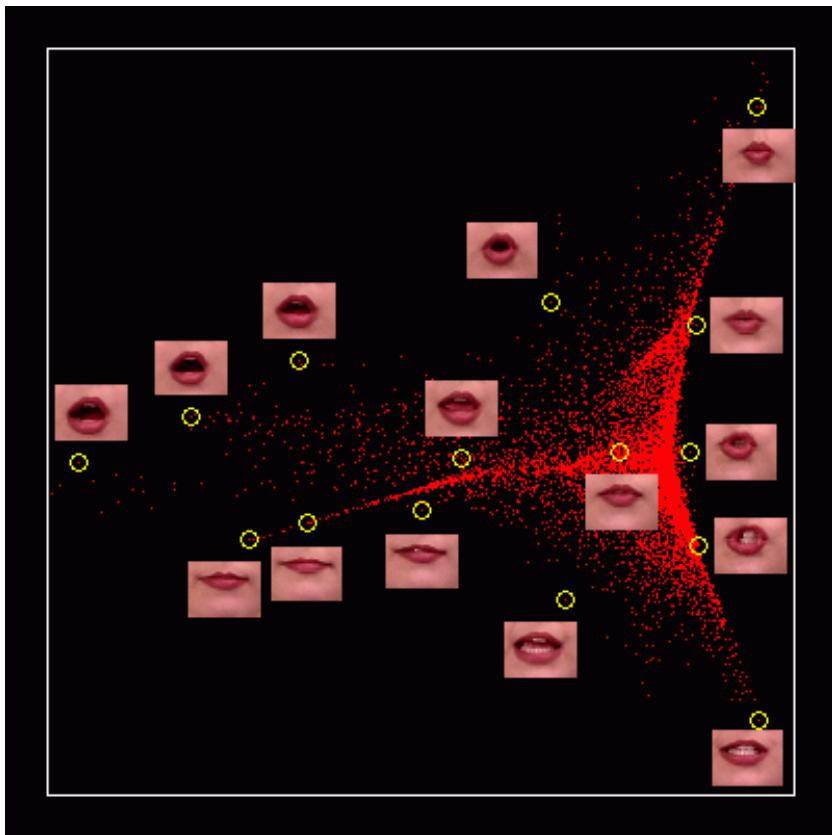
$$J_1(\mathbf{W}) = \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{j=1}^K W_{ji} \mathbf{x}_{j(i)} \right\|^2$$

Step 3. Get $\{y_1, y_2, \dots, y_n\} \in \mathbb{R}^q$
by minimizing

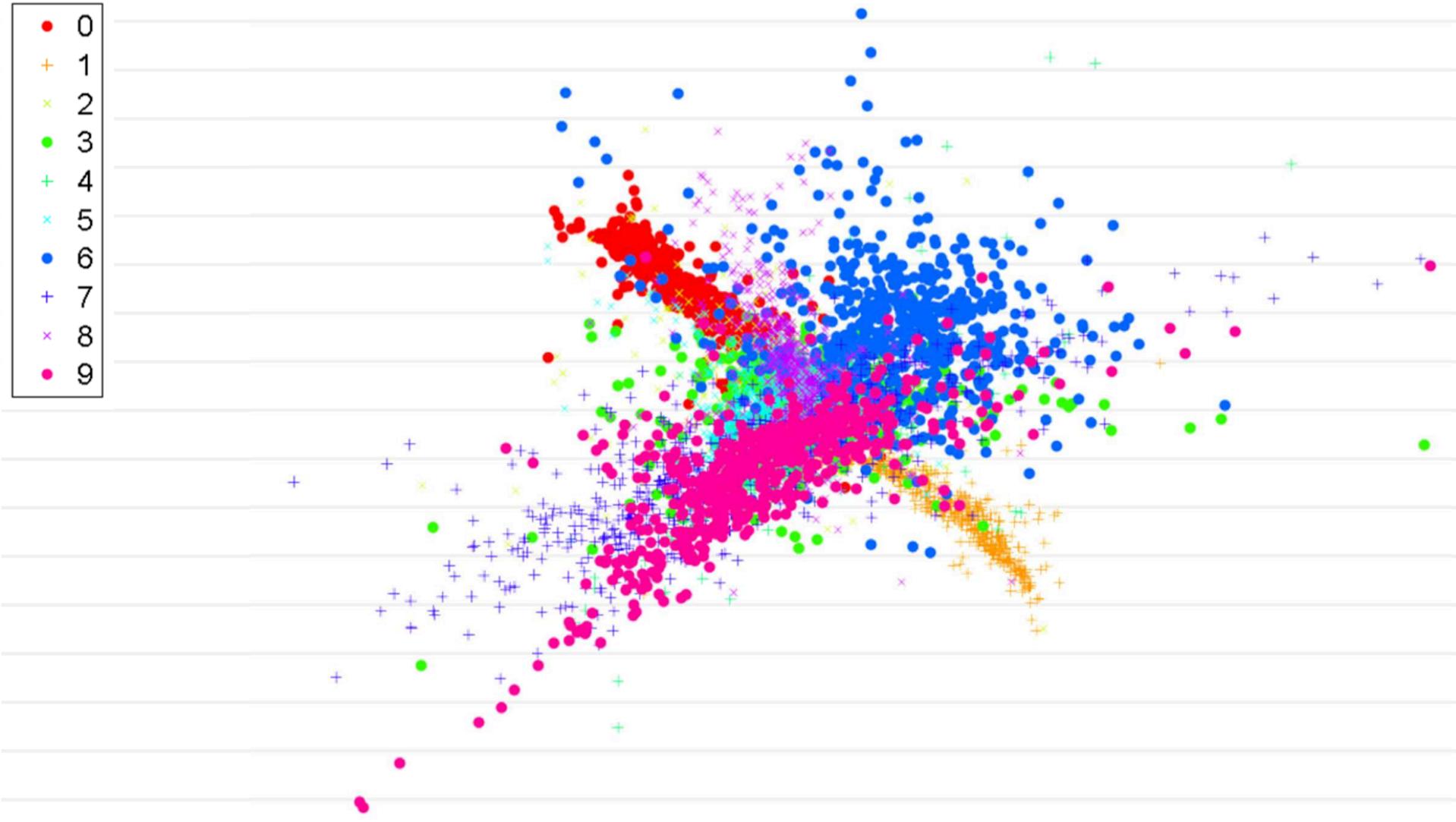
$$J_2(\mathbf{Y}) = \sum_{i=1}^n \left\| \mathbf{y}_i - \sum_{j=1}^N W_{ji} \mathbf{y}_{j(i)} \right\|^2$$

with standard normalizing conditions

LLE: Lips pictures



LLE: MNIST



LLE: pro and con

- no local minima, one free parameter
- incremental & fast
- simple linear algebra operations
- can distort global structure

Laplacian Eigenmaps (LE) (M. Belkin, P. Niyogi, 2002)

Step 1 is standard

Step 2: $w_{ij} = w_{0,ij} \times \exp(-t^{-1} \times \|X_i - X_j\|^2)$ for all $(X_i, X_j) \in V$

Step LE3. Construct embedding. Values $y_i = h(X_i)$ minimize quadratic form

$$\sum_{i,j=1}^n W_{ij} \times \|y_i - y_j\|^2 = \text{Tr}(Y^T \times R_{LE} \times Y)$$

s.t. $Y^T \times D \times Y = I_q$

$D = \text{Diag}(\sum_{j=1}^n W_{ij})$, where $R_{LE} = D - W$, $W = \|w_{ij}\|$,

Y^T is $(q \times n)$ -matrix with columns $\{y_1, y_2, \dots, y_n\} \in R^q$ 59

Solution: $\mathbf{Y} = (Y_1, Y_2, \dots, Y_q)$, $Y_i \in \mathbb{R}^n$ is a solution of a system of n linear equations

$$R_{LE} \mathbf{Y} = \lambda_i \times D\mathbf{Y}$$

for q smallest **positive** roots $\lambda_1, \lambda_2, \dots, \lambda_q > 0$ of the equation

$$\text{Det}(R_{LE} - \lambda \times D) = 0.$$

We do not consider $\lambda_0 = 0$, since it corresponds to $Y_0 = \mathbf{1} \in \mathbb{R}^n$ and contradicts our normalization $\mathbf{Y} \times \mathbf{1} = \mathbf{0}$

Justification: why Laplacian

- Manifold $\mathbf{X} \subset \mathbb{R}^p$ is smooth with Riemannian structure
- Let $q = 1$ and h is an embedding transformation of \mathbf{X} in $\mathbf{Y} = h(\mathbf{X}) \subset \mathbb{R}^1$,
 $\nabla h(x)$ is a gradient (vector field on a manifold \mathbf{X}):

$$\|h(x) - h(z)\| \leq \|\nabla h(x)\| \times \|x - z\| + o\|x - z\|;$$

- “we see that if $\|\nabla h(x)\|$ provides us with an estimate of how far apart h maps nearby points. We therefore **look a map that best preserves locality on average by trying to find**

$$\arg \min \{I(h) : \|h\|_{L^2(\mathbf{X})} = 1\} \text{ , where } I(h) = \int_{\mathbf{X}} \|\nabla h(x)\|^2 mes(dx);$$

- $L_{LB}h = -\operatorname{div}(\nabla h)$, then by Stocks theorem

$$\int_X \|\nabla h\|^2 = \int_X L_{LB}(h) \times h$$

\Rightarrow minimization of $I(h) \Leftrightarrow$ minimization of $\int_X L_{LB}(h) \times h$

- operator L_{LB} has discrete spectra $0 = \lambda_0 \leq \lambda_1 \leq \lambda_2 \leq \dots$, and $h_0 = \text{const}$
 $\Rightarrow h_1 \perp h_0$ minimizes $I(h)$, h_i are eigenfunctions of L_{LB} ;

- **minimization of $I(h) = \int_X \|\nabla h\|^2 \Leftrightarrow$ minimization of a sample operator**

$$Lh = \sum_{i,j=1}^n W_{ij} \times (h(X_i) - h(X_j))^2$$

defined on nodes of a graph графа $\Gamma(X_n)$; it is a discrete analogue of Laplace-Beltrami operator

Hessian Eigenmaps, HE (D.L. Donoho, C. Grimes, 2002).

Hessian Eigenmaps is “LE only with the Hessian replacing the Laplacian”

$$\int_X \|\nabla h(x)\|^2 mes(dx) \Rightarrow H(h) = \int_X \|H_h(x)\|_F^2 mes(dx)$$

Two more methods (not covered here):

Riemannian Manifold Learning, RML

Local Tangent Space Alignment, LTSA

stochastic neighbor embedding

“Stochastic” similarity in the initial space

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

Similarity in the compressed space

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

stochastic neighbor embedding

“Stochastic” similarity in the initial space

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

Similarity in the compressed space

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

$$Cost = \sum_i KL(P_i \| Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

stochastic neighbor embedding

$$\frac{\partial Cost}{\partial y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j)$$

To escape poor minima:

- Initialize with Gaussian noise
- Add diminishing Gaussian noise on each stage

t-sne: Similarity in the compressed space in case of t(1) distribution

$$q_{ij} = \frac{(1 + \| y_i - y_j \|^2)^{-1}}{\sum_{k \neq l} (1 + \| y_k - y_l \|^2)^{-1}}$$

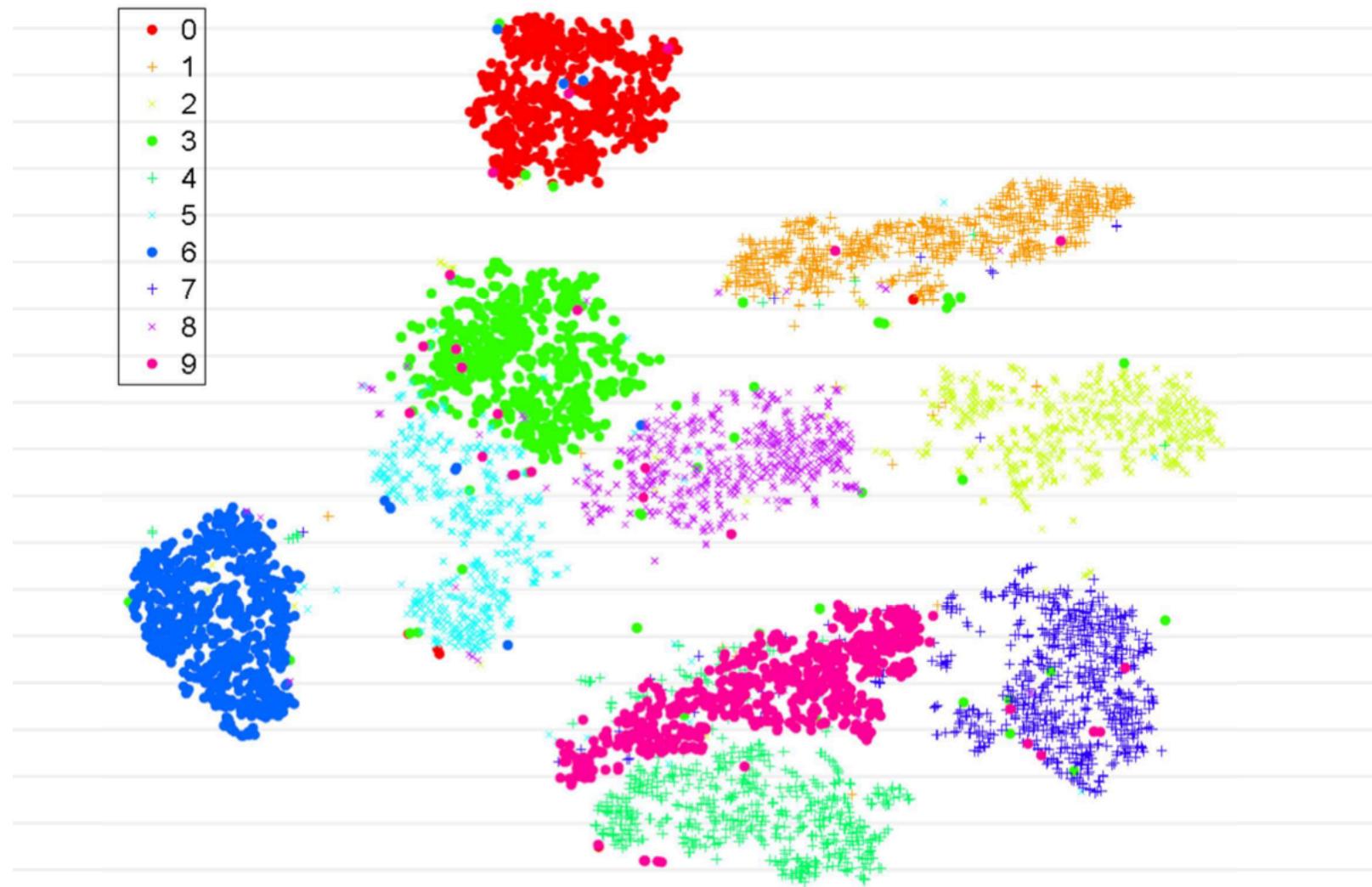
stochastic neighbor embedding

t-sne: minimize KL between the whole densities P and Q

$$Cost = KL(P \parallel Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$$

t-stochastic neighbor embedding: MNIST



No Free Lunch

the “curvier” your manifold, the denser your data must be

