

Nonlinear Dimensionality Reduction and Euclidean Embedding

Outline

- Euclidean Embedding
- Dimensionality Reduction vs Euclidean Embedding
- Multidimensional Scaling
- Kernel PCA
- ISOMAP
- Laplacian Eigenmaps

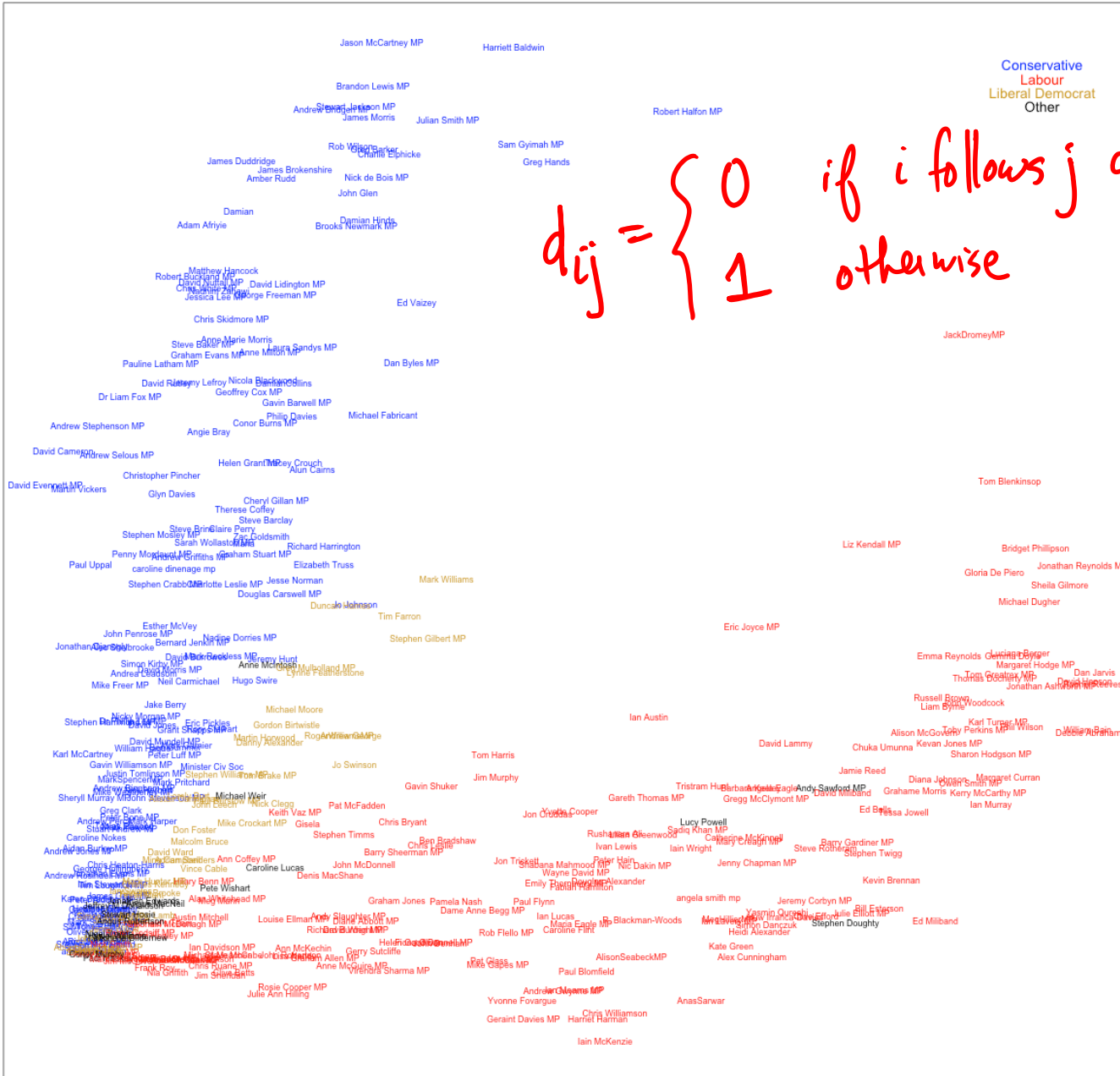
Euclidean Embedding

- A dissimilarity matrix is a square, $n \times n$ matrix $\mathbf{D} = [d_{ij}]$ such that
 - $d_{ij} \geq 0$
 - $d_{ij} = d_{ji}$
 - $d_{ii} = 0$
- Given: n objects (not necessarily Euclidean vectors) with dissimilarity matrix \mathbf{D} , find Euclidean vectors $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^p$ such that

$$d_{ij} \approx \|\mathbf{y}_i - \mathbf{y}_j\|_2$$

- Uses?
 - 1) visualization / exploratory data analysis ($p=2$ or 3)
 - 2) apply an algorithm that works on Euclidean data

Two dimensional clustering of UK Members of Parliament



Based on Twitter follows by UK Members of Parliament, Aug 2013. <http://andrewwhitby.com/2013/08/21/clustering-of-uk-mps/>

- Two dimensional Euclidean embedding of UK members of Parliament
- Dissimilarity based on Twitter follows
- Algorithm is classical MDS
- Color indicates party affiliation (not used by algorithm)

Euclidean Embedding vs Dimensionality Reduction

- Differences:

- DR starts with a high dimensional data matrix $X \in \mathbb{R}^{d \times n}$
- EE: Starts with a dissimilarity matrix
- EE methods apply to non Euclidean data

- Similarities:

Same goal: low dimensional representation of data

- An EE method can be used for DR, but not vice versa

Euclidean Distance Matrices

- An $n \times n$ dissimilarity matrix \mathbf{D} is called a *Euclidean distance matrix* if there exists p and $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^p$ such that

$$d_{ij} = \|\mathbf{y}_i - \mathbf{y}_j\|_2$$

for all i, j .

- **Theorem (Part 1):** Let \mathbf{D} be an $n \times n$ dissimilarity matrix. Set $\mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H}$ where

$$\mathbf{A} = [a_{ij}], \quad a_{ij} = -d_{ij}^2, \quad \mathbf{H} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T.$$

Then \mathbf{D} is a Euclidean distance matrix iff \mathbf{B} is PSD.

Euclidean Distance Matrices

- **Theorem (Part 2):** If B is PSD with positive eigenvalues $\lambda_1 > \dots > \lambda_p$ and corresponding eigenvectors

$$\mathbf{u}_1 = \begin{bmatrix} u_{11} \\ \vdots \\ u_{1n} \end{bmatrix}, \dots, \mathbf{u}_p = \begin{bmatrix} u_{p1} \\ \vdots \\ u_{pn} \end{bmatrix}$$

normalized such that

$$\mathbf{u}_k^T \mathbf{u}_k = \lambda_k$$

then the vectors

$$\mathbf{y}_i = (u_{1i}, \dots, u_{pi})^T \in \mathbb{R}^p$$

satisfy

$$d_{ij} = \|\mathbf{y}_i - \mathbf{y}_j\|$$

- Proof: Mardia, Kent, and Bibby, *Multivariate Analysis*, 1979.

Classical MDS

- Even if \mathbf{D} is not a Euclidean distance matrix, the previous result suggests an algorithm for MDS:
- Input: \mathbf{D}, p
 - Form \mathbf{B} as in the theorem
 - Compute the eigenvalue decomposition

$$\mathbf{B} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T, \quad (n \times n)$$

where $\mathbf{V} = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_n]$, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$

- Set $\mathbf{u}_k = \sqrt{\lambda_k} \mathbf{v}_k$

Return $\mathbf{y}_i = i^{\text{th}}$ row of

$$\mathbf{U}_p = [\mathbf{u}_1 \ \cdots \ \mathbf{u}_p] \quad (n \times p).$$

- True (A) or False (B): The above procedure is valid for any dissimilarity matrix \mathbf{D}

Stress-Based MDS

- A second approach to MDS is to minimize the *stress* objective function

$$\min_{y_1, \dots, y_n \in \mathbb{R}^p} \sum_{1 \leq i < j \leq n} w_{ij} (d_{ij} - \|y_i - y_j\|)^2$$

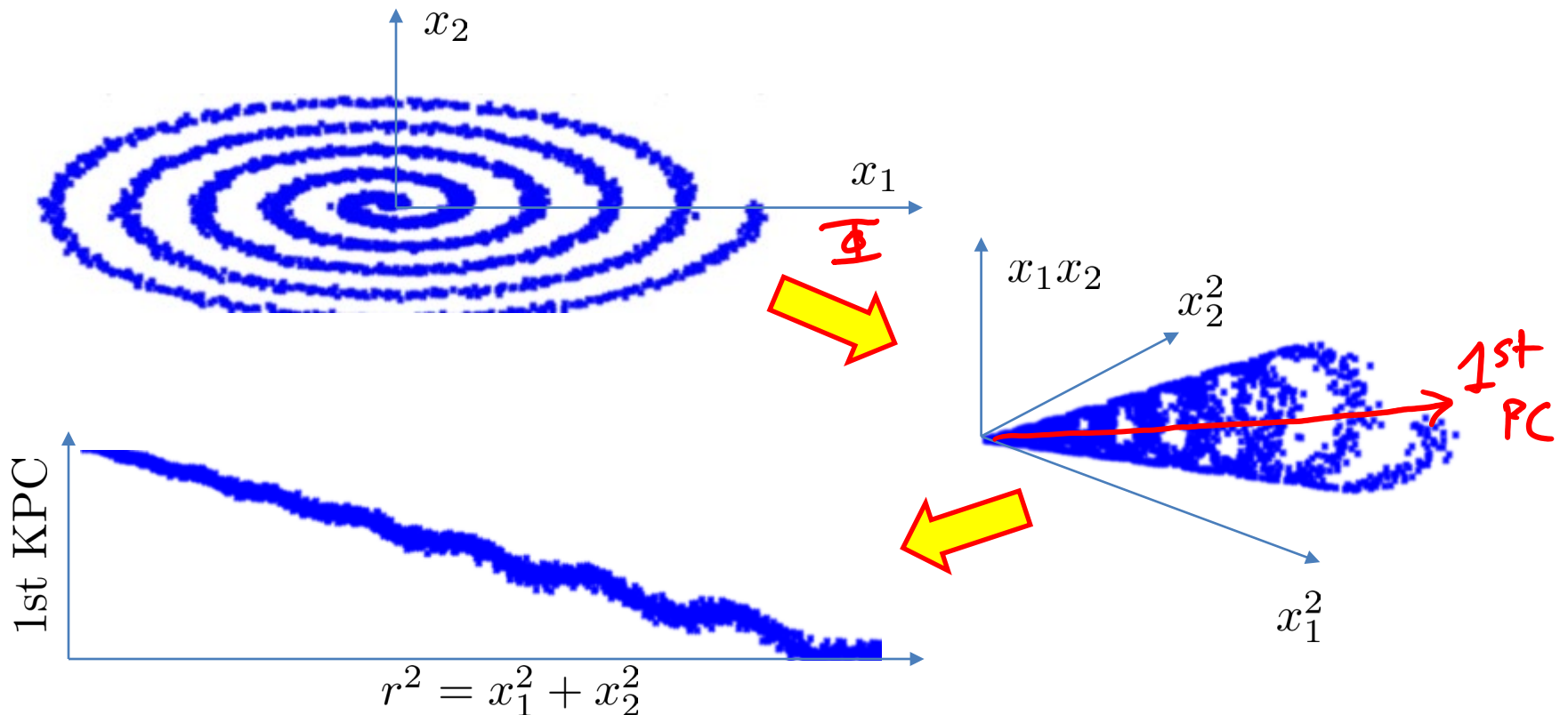
where $w_{ij} \geq 0$ are weights reflecting the similarity of the data points

- Examples: $w_{ij} = 1$ or $w_{ij} = d_{ij}^{-\alpha}$, $\alpha > 0$.
- This problem is nonconvex, but a local minimizer can be obtained efficiently with a majorize-minimize algorithm.

Demo

Kernel PCA

- Intuitively, you can think of this like any other kernel method: First apply a nonlinear feature map Φ (associated to a SPD kernel), and then do PCA in the new feature space.
- KPCA is a good general purpose method for dimensionality reduction



KPCA Algorithm

Input: $\mathbf{x}_1, \dots, \mathbf{x}_n$, dimension p , **kernel k**

$$[K]_{ij} = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$$

- Form centered kernel matrix

$$\tilde{K} = K - OK - KO + OKO$$

$$[\tilde{K}]_{ij} = \langle \tilde{\Phi}(\mathbf{x}_i), \tilde{\Phi}(\mathbf{x}_j) \rangle$$

where O is $n \times n$ and has entries equal to $\frac{1}{n}$

- Compute the eigenvalue decomposition

where

$$\tilde{\Phi}(\mathbf{x}) = \Phi(\mathbf{x}) - \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$$

$$\tilde{K} = U \Lambda U^T$$

$$U = [\mathbf{u}_1 \quad \dots \quad \mathbf{u}_n]$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m, 0, \dots, 0)$$

- Set $\boldsymbol{\alpha}_j = (\alpha_{j1}, \dots, \alpha_{jn})^T := \frac{1}{\sqrt{\lambda_j}} \mathbf{u}_j \in \mathbb{R}^n$, $1 \leq j \leq p$

Output: Dimensionality reduction mapping

$$\mathbf{x} \mapsto \mathbf{y} = (y_1, \dots, y_p)^T \in \mathbb{R}^p$$

where

$$y_j = \sum_{i=1}^n \alpha_{ji} \tilde{k}(\mathbf{x}, \mathbf{x}_i)$$

$$\tilde{k}(\mathbf{x}, \mathbf{x}') = \langle \tilde{\Phi}(\mathbf{x}), \tilde{\Phi}(\mathbf{x}') \rangle$$

KPCA Algorithm

Input: $\mathbf{x}_1, \dots, \mathbf{x}_n$, dimension p

- Form centered kernel matrix

$$\tilde{K} = K - OK - KO + OKO$$

where O is $n \times n$ and has entries equal to $\frac{1}{n}$

- Compute the eigenvalue decomposition

$$\tilde{K} = U\Lambda U^T$$

$$U = [\mathbf{u}_1 \quad \cdots \quad \mathbf{u}_n]$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m, 0, \dots, 0)$$

- Set $\boldsymbol{\alpha}_j = (\alpha_{j1}, \dots, \alpha_{jn})^T := \frac{1}{\sqrt{\lambda_j}} \mathbf{u}_j \in \mathbb{R}^n$, $1 \leq j \leq p$

Output: Dimensionality reduction mapping

$$\mathbf{x} \mapsto \mathbf{y} = (y_1, \dots, y_p)^T \in \mathbb{R}^p$$

where

$$y_j = \sum_{i=1}^n \alpha_{ji} \tilde{k}(\mathbf{x}, \mathbf{x}_i)$$

True (A) or False (B):
This procedure is valid as long as $p \leq m$, the number of nonzero eigenvalues of the centered kernel matrix.

Isomap

- KPCA does not try to capture the *intrinsic dimension* of the data
- A method that does is Isomap (isometric mapping)
- The algorithm is simple:
 - Input $\mathbf{x}_1, \dots, \mathbf{x}_n$
 - Construct a similarity graph such as a k -nearest neighbor graph
 - Form a dissimilarity matrix $\mathbf{D} = [d_{ij}]$ where d_{ij} = length of shortest path connecting \mathbf{x}_i and \mathbf{x}_j
 - Apply MDS to \mathbf{D} to get an embedding $\mathbf{y}, \dots, \mathbf{y}_n$
- Tenenbaum, de Silva and Langford, “A global geometric framework for nonlinear dimensionality reduction,” *Science*, 290, 2319-2323 (2000)

Isomap

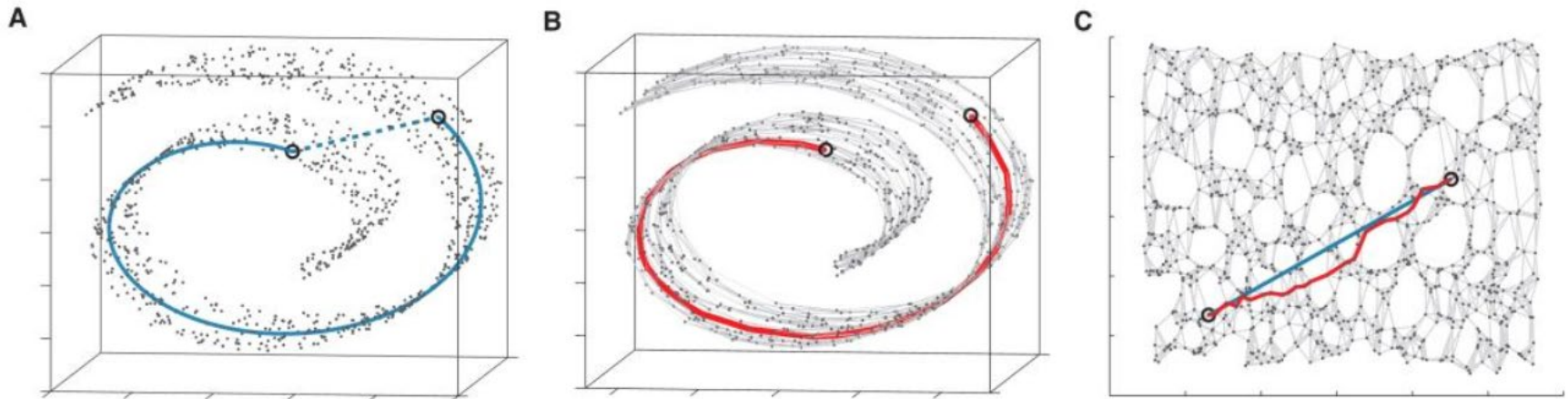


Fig. 3. The "Swiss roll" data set, illustrating how Isomap exploits geodesic paths for nonlinear dimensionality reduction. **(A)** For two arbitrary points (circled) on a nonlinear manifold, their Euclidean distance in the high-dimensional input space (length of dashed line) may not accurately reflect their intrinsic similarity, as measured by geodesic distance along the low-dimensional manifold (length of solid curve). **(B)** The neighborhood graph G constructed in step one of Isomap (with $K = 7$ and $N =$

1000 data points) allows an approximation (red segments) to the true geodesic path to be computed efficiently in step two, as the shortest path in G . **(C)** The two-dimensional embedding recovered by Isomap in step three, which best preserves the shortest path distances in the neighborhood graph (overlaid). Straight lines in the embedding (blue) now represent simpler and cleaner approximations to the true geodesic paths than do the corresponding graph paths (red).

Isomap

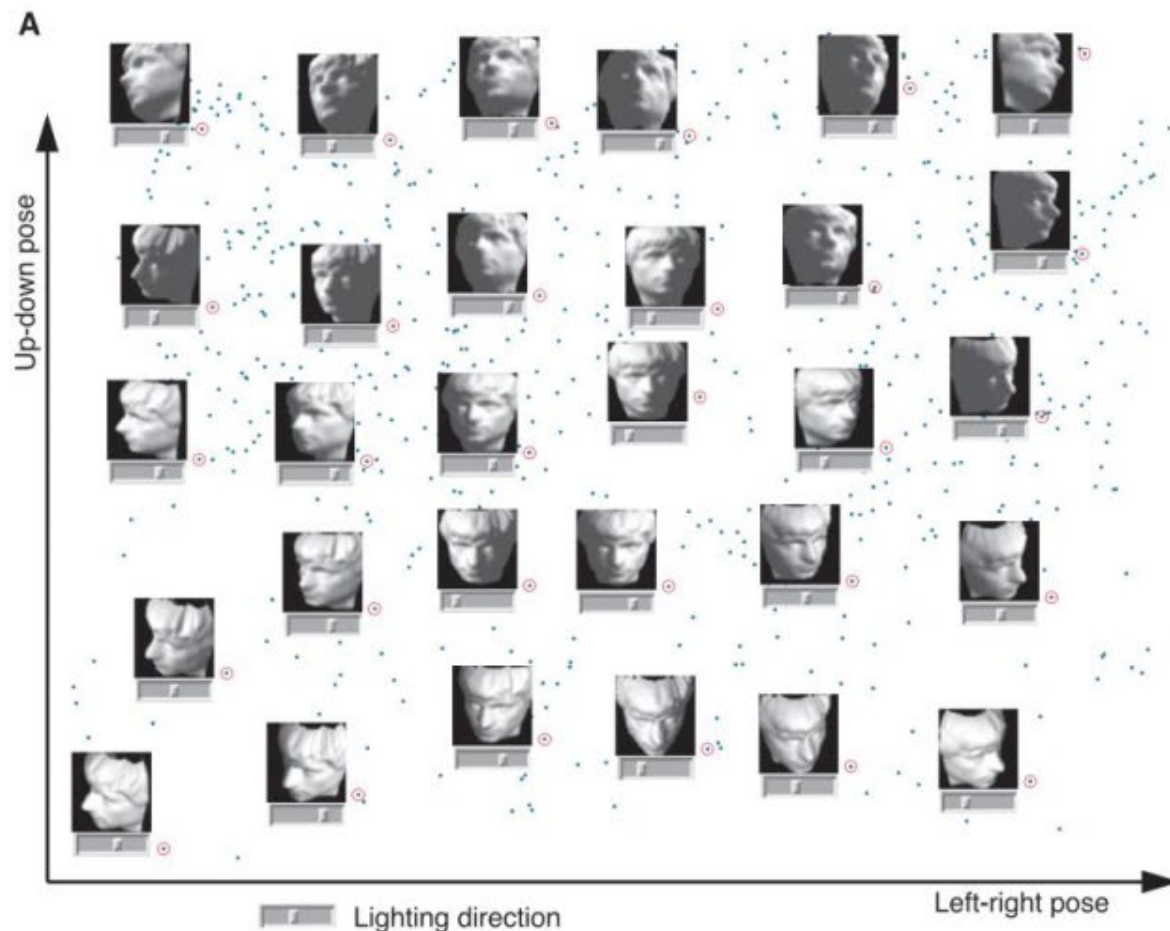


Fig. 1. (A) A canonical dimensionality reduction problem from visual perception. The input consists of a sequence of 4096-dimensional vectors, representing the brightness values of 64 pixel by 64 pixel images of a face rendered with different poses and lighting directions. Applied to $N = 698$ raw images, Isomap ($K = 6$) learns a three-dimensional embedding of the data's intrinsic geometric structure. A two-dimensional projection is shown, with a sample of the original input images (red circles) superimposed on all the data points (blue) and horizontal sliders (under the images) representing the third dimension. Each coordinate axis of the embedding correlates highly with one degree of freedom underlying the original data: left-right pose (x axis, $R = 0.99$), up-down pose (y axis, $R = 0.90$), and lighting direction (slider position, $R = 0.92$). The input-space distances $d_x(i, j)$ given to Isomap were Euclidean distances between the 4096-dimensional image vectors.

Laplacian Eigenmaps

- This is the Euclidean embedding method implicit in spectral clustering
- Discussed in spectral clustering notes
- Can be derived independent from spectral clustering, as follows:
- Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be training data, not necessarily Euclidean vectors
- Let $\mathbf{W} = [w_{ij}]$ be a similarity graph as in spectral clustering
- Idea: select $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^p$ to optimize

$$\min_{\mathbf{y}_1, \dots, \mathbf{y}_n} \frac{1}{2} \sum_{i,j=1}^n w_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|^2 = \text{tr}(\mathbf{Y} \mathbf{L} \mathbf{Y}^T)$$
$$\mathbf{Y} = [\mathbf{y}_1 \cdots \mathbf{y}_n]$$

- To avoid a trivial solution, need to impose an energy constraint. Options:

$$\mathbf{Y} \mathbf{Y}^T = \mathbf{I}_p, \quad \mathbf{Y} \mathbf{D} \mathbf{Y}^T = \mathbf{I}_p$$

Laplacian Eigenmaps

- The solution to

$$\begin{aligned} \min_{\mathbf{Y}} \quad & \text{tr}(\mathbf{Y} \mathbf{L} \mathbf{Y}^T) \\ \text{s.t.} \quad & \mathbf{Y} \mathbf{Y}^T = \mathbf{I} \end{aligned}$$

is

where $\mathbf{u}_1, \dots, \mathbf{u}_p$ are the p smallest eigenvectors of

- If the graph is connected, then \mathbf{u}_1 is a multiple of $\mathbf{1}$ and we can equivalently solve

which determines a dimensional embedding

Laplacian Eigenmaps

- The solution to

$$\begin{aligned} \min_{\mathbf{Y}} \quad & \text{tr}(\mathbf{Y} \mathbf{L} \mathbf{Y}^T) \\ \text{s.t.} \quad & \mathbf{Y} \mathbf{D} \mathbf{Y}^T = \mathbf{I} \end{aligned}$$

is

where $\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_p$ are the p smallest eigenvectors of

- If the graph is connected, then \mathbf{u}_1 is a multiple of $\mathbf{1}$ and we can equivalently solve

which again determines a $p - 1$ dimensional embedding

Laplacian Eigenmaps

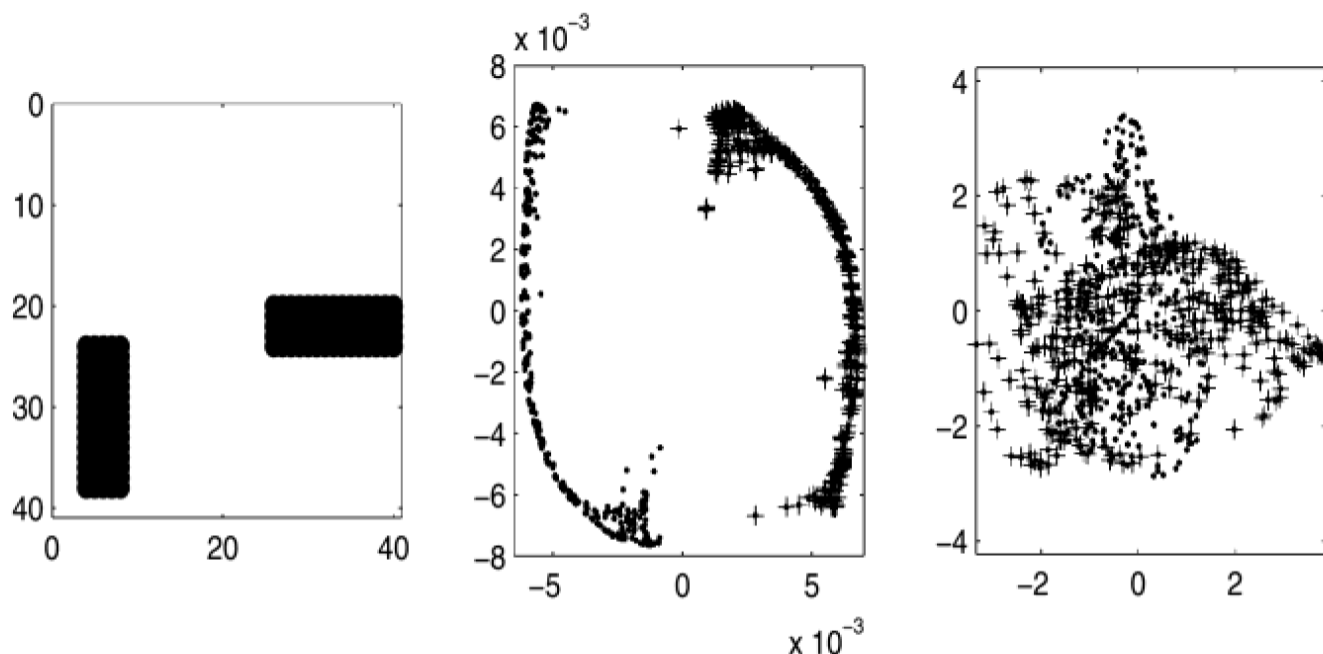


Figure 3: (Left) A horizontal and a vertical bar. (Middle) A two-dimensional representation of the set of all images using the Laplacian eigenmaps. (Right) The result of PCA using the first two principal directions to represent the data. Blue dots correspond to images of vertical bars, and plus signs correspond to images of horizontal bars.

t-SNE

- Given $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- For all $i \neq j$, compute

$$p_{i|j} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$$

- Set $p_{ij} = (p_{i|j} + p_{j|i}) / 2n$. These values are nonnegative, symmetric, and sum to 1 (summing over both i and j) \rightarrow probability distribution P
- For candidate embedded points $\mathbf{y}_1, \dots, \mathbf{y}_n$, define a probability distribution Q analogously but using

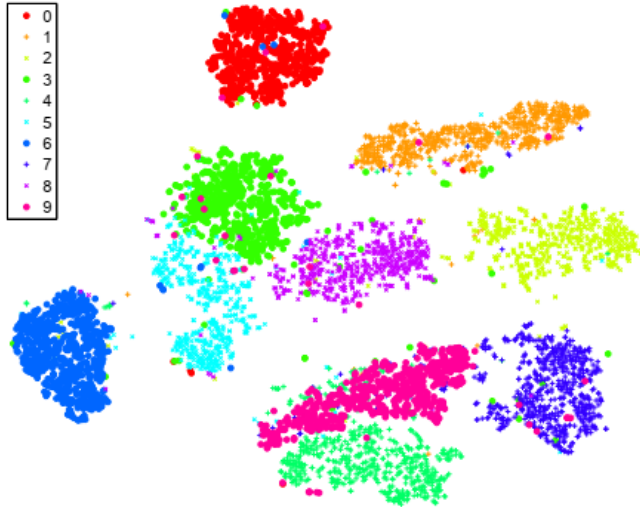
$$q_{i|j} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|\mathbf{y}_i - \mathbf{y}_k\|^2)^{-1}}$$

- Optimize $\mathbf{y}_1, \dots, \mathbf{y}_n$ by minimizing the Kullback-Liebler divergence, $KL(P\|Q) = \sum_{ij} p_{ij} \log(p_{ij}/q_{ij})$ using gradient descent
- Remarks: data-dependent bandwidths, name, avoid crowding

t-SNE avoids
crowding

t-SNE Illustration

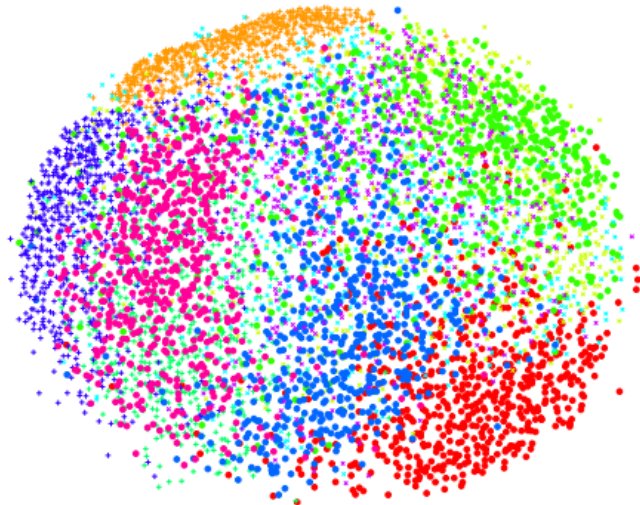
Van der Maaten
and Hinton, 2008



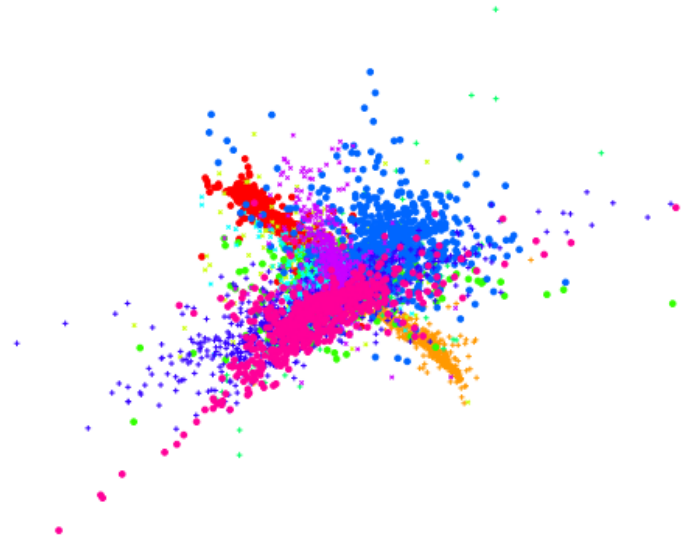
(a) Visualization by t-SNE.



(a) Visualization by Isomap.



(b) Visualization by Sammon mapping.



t-SNE Demo

<https://distill.pub/2016/misread-tsne/>

Summary

- DR and EE: Two classes of methods for learning a low-dimensional Euclidean representation of data
- Methods differ in whether they emphasize preservation of local or global distances
- Some methods have natural out-of-sample extensions, others do not
- Some methods aim to capture intrinsic dimension (manifold learning), some do not
- t-SNE and UMAP (not covered) are quite popular. These use gradient descent to iteratively optimize some nonconvex objective. Initialization is important, e.g., by LEM, which globally optimizes its objective.
 - https://umap-learn.readthedocs.io/en/latest/how_umap_works.html
 - <https://arxiv.org/abs/1802.03426>
 - <https://arxiv.org/abs/2009.12981>

Autoencoder