

Lecture 7: Nonlinear Dimensionality Reduction

Shukai Gong

1 Manifold Learning

1.1 Multi-dimensional Scaling

Metric MDS

Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$, we can compute a distance matrix

$$\mathbf{D} = [d_{ij}] \in \mathbb{R}^{N \times N}, \quad d_{ij} = d(\mathbf{x}_i, \mathbf{x}_j)$$

Metric MDS aims at finding low-dimensional latent representation $\{\mathbf{z}_n\}_{n=1}^N$ to keep isometry as much as possible via

$$\min_{\{\mathbf{z}_n\}_{n=1}^N} \text{Stress}_d(\{\mathbf{z}_n\}_{n=1}^N) = \left(\min_{\{\mathbf{z}_n\}_{n=1}^N} \sum_{i \neq j} (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|_p)^2 \right)^{\frac{1}{2}}$$

where $p = 1, 2$ in general.

[Note]

- There's no explicit expression for $\{\mathbf{x}_n\} \rightarrow \{\mathbf{z}_n\}$
- There isn't unique solution for $\{\mathbf{z}_n\}$. For example, if we take $p = 2$ and

$$\{\mathbf{z}_n^*\} = \arg \min_{\{\mathbf{z}_n\}_{n=1}^N} \left(\sum_{i \neq j} (d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|_2)^2 \right)^{\frac{1}{2}}$$

and U as any unitary matrix ($U^*U = \mathbf{I}$), then $\{U\mathbf{z}_n^*\}$ is also a solution since $\|U\mathbf{z}_i - U\mathbf{z}_j\|_2 = \|\mathbf{z}_i - \mathbf{z}_j\|_2$.

Classic MDS

Classic MDS is a special case of Metric MDS where $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$ is Euclidean. We replace our optimization goal from $\min \text{Stress}_d(\{\mathbf{z}_n\}_{n=1}^N)$, **which minimizes the difference between pairwise distances in the original space and the latent space**, to

$$\min \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) = \min_{\{\mathbf{z}_n\}_{n=1}^N} \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^\top \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{\frac{1}{2}}$$

which minimizes the difference between inner product in the original space and the latent space.

Denote our dataset as $\mathbf{X} \in \mathbb{R}^{N \times D}$. Here the Gram Matrix is defined as $\mathbf{K} = [k_{ij}] = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C}$ with **centering matrix** $\mathbf{C} = \mathbf{I}_N - \frac{1}{N}\mathbf{1}_{N \times N}$. The low-dimension embedding \mathbf{Z}^* is derived first by performing EVD on $\mathbf{K} := \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$, then

$$\mathbf{Z}^* = \mathbf{V}_L \mathbf{\Lambda}_L^{\frac{1}{2}}$$

Denote $\tilde{\mathbf{X}} = \mathbf{C}\mathbf{X}$, then $\mathbf{K} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$ (See Appendix for derivation). Back to our optimization goal of

$$\begin{aligned} \min_{\{\mathbf{z}_n\}_{n=1}^N} \text{Strain}_d(\{\mathbf{z}_n\}_{n=1}^N) &= \min_{\{\mathbf{z}_n\}_{n=1}^N} \left(\frac{\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^\top \mathbf{z}_j)^2}{\sum_{i,j=1}^N k_{ij}^2} \right)^{\frac{1}{2}} = \min_{\{\mathbf{z}_n\}_{n=1}^N} \left(\sum_{i,j=1}^N (k_{ij} - \mathbf{z}_i^\top \mathbf{z}_j)^2 \right)^{\frac{1}{2}} \\ &= \min_{\mathbf{Z}} \|\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top\|_F = \min_{\mathbf{Z}} \|\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top\|_F^2 \\ &= \min_{\mathbf{Z}} \text{tr}[(\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top)^\top (\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top)] = \min_{\mathbf{Z}} \text{tr}[(\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top)^2] \end{aligned}$$

Performing EVD on \mathbf{K} and $\mathbf{Z}\mathbf{Z}^\top$, we have

$$\mathbf{K} = \mathbf{V}\mathbf{\Delta}\mathbf{V}^\top, \mathbf{Z}\mathbf{Z}^\top = \mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top$$

and then

$$\begin{aligned} \|\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top\|_F^2 &= \text{tr}[(\mathbf{V}\mathbf{\Delta}\mathbf{V}^\top - \mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top)^2] = \text{tr}[(\mathbf{V}\mathbf{\Delta}\mathbf{V}^\top - \mathbf{V}\mathbf{V}^\top\mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top\mathbf{V}\mathbf{V}^\top)^2] \\ &= \text{tr}[(\mathbf{V}(\mathbf{\Delta} - \mathbf{V}^\top\mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top\mathbf{V})\mathbf{V}^\top)^2] = \text{tr}[\mathbf{V}^2(\mathbf{\Delta} - \mathbf{V}^\top\mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top\mathbf{V})^2(\mathbf{V}^\top)^2] \\ &= \text{tr}[(\mathbf{V}^\top)^2\mathbf{V}^2(\mathbf{\Delta} - \mathbf{V}^\top\mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top\mathbf{V})^2] = \text{tr}[(\mathbf{\Delta} - \mathbf{V}^\top\mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top\mathbf{V})^2] \end{aligned}$$

Let $\mathbf{M} := \mathbf{V}^\top\mathbf{Q}$, then

$$\begin{aligned} \min_{\mathbf{Z}} \|\mathbf{K} - \mathbf{Z}\mathbf{Z}^\top\|_F^2 &= \min_{\mathbf{M}, \mathbf{\Psi}} \text{tr}[(\mathbf{\Delta} - \mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)^2] \\ &= \min_{\mathbf{M}, \mathbf{\Psi}} \text{tr}(\mathbf{\Delta}^2) - 2\text{tr}(\mathbf{\Delta}\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top) + \text{tr}[(\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)^2] \end{aligned}$$

Denote $\mathcal{L} = \text{tr}(\mathbf{\Delta}^2) - 2\text{tr}(\mathbf{\Delta}\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top) + \text{tr}[(\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)^2]$. First we take the derivative w.r.t. \mathbf{M} and set it to zero:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{M}} &= -2\mathbf{\Delta}\mathbf{M}\mathbf{\Psi} + 2(\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)\mathbf{M}\mathbf{\Psi} = 0 \\ \Rightarrow \mathbf{M}\mathbf{\Psi}\mathbf{M}^\top &= \mathbf{\Delta} \end{aligned}$$

Before taking the derivative w.r.t. $\mathbf{\Psi}$, we first change \mathcal{L} into:

$$\begin{aligned} \mathcal{L} &= \text{tr}(\mathbf{\Delta}^2) - 2\text{tr}(\mathbf{\Delta}\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top) + \text{tr}[(\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)^2] \\ &= \text{tr}(\mathbf{\Delta}^2) - 2\text{tr}(\mathbf{M}^\top\mathbf{\Delta}\mathbf{M}\mathbf{\Psi}) + \text{tr}[(\mathbf{M}^\top\mathbf{M}\mathbf{\Psi})^2] \end{aligned}$$

then

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{\Psi}} &= -2\mathbf{M}^\top\mathbf{\Delta}\mathbf{M} + 2(\mathbf{M}^\top\mathbf{M}\mathbf{\Psi})\mathbf{M}^\top\mathbf{M} \\ &= -2\mathbf{M}^\top\mathbf{\Delta}\mathbf{M} + 2\mathbf{M}^\top(\mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)\mathbf{M} = 0 \\ \Rightarrow \mathbf{M}^\top\mathbf{\Psi}\mathbf{M} &= \mathbf{\Delta} \end{aligned}$$

Both FOC points to $\mathbf{M}^\top\mathbf{\Psi}\mathbf{M} = \mathbf{\Delta}$. One possible solution to this is

$$\mathbf{M} = \mathbf{I}, \mathbf{\Psi} = \mathbf{\Delta}$$

which means that the minimum of the non-negative objective function $\text{tr}[(\mathbf{\Delta} - \mathbf{M}\mathbf{\Psi}\mathbf{M}^\top)^2]$ is 0. Therefore, we have

$$\mathbf{M} = \mathbf{I} = \mathbf{V}^\top\mathbf{Q} \Rightarrow \mathbf{Q} = \mathbf{V}$$

Recall that

$$\mathbf{Z}\mathbf{Z}^\top = \mathbf{Q}\mathbf{\Psi}\mathbf{Q}^\top = \mathbf{V}\mathbf{\Delta}\mathbf{V}^\top = \mathbf{V}\mathbf{\Delta}^{\frac{1}{2}}\mathbf{\Delta}^{\frac{1}{2}}\mathbf{V}^\top \Rightarrow \mathbf{Z} = \mathbf{V}\mathbf{\Delta}^{\frac{1}{2}}$$

Truncating this \mathbf{Z} gives us $\mathbf{Z}^* = \mathbf{V}_L\mathbf{\Delta}_L^{\frac{1}{2}} \in \mathbb{R}^{N \times L}$.

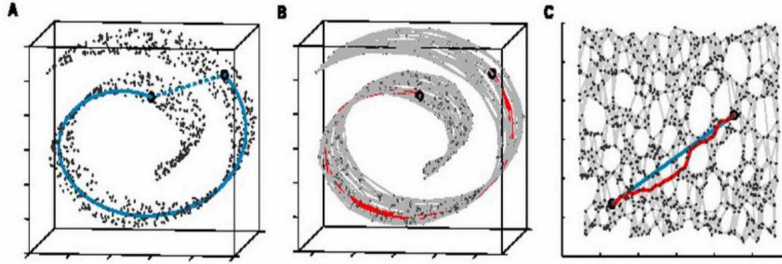
1.2 ISOMAP

ISOMAP

ISOMAP is a special case of MDS where **isometry is kept under geodesic distance** as much as possible. Given a set of data $\{\mathbf{x}_n\}_{n=1}^N$

1. Determine the neighbors of each data point and construct a K -nearest neighbor (KNN) graph of the data.
2. Compute the shortest path (**Dijkstra/Floyd**) distance between arbitrary two nodes and obtain an approximate geodesic distance matrix $\mathbf{D} = [d_{ij}] \in \mathbb{R}^{N \times N}$.
3. Compute low-dimensional embedding by MDS similarly

$$\begin{cases} \mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C} \\ \mathbf{K} = \mathbf{V}\mathbf{\Delta}\mathbf{V}^\top \end{cases} \Rightarrow \mathbf{Z}^* = \mathbf{V}_L \mathbf{\Delta}_L^{\frac{1}{2}}$$



ISOMAP

1.3 Locally Linear Embedding

Locally Linear Embedding (LLE)

LLE keeps isometry indirectly through inheriting **local linear self-representation power**. Local linear self-representation means that each data point can be represented by a linear combination of its neighbors: given a sample \mathbf{x}_i and its K neighbors $\mathbf{X}_i = [\mathbf{x}_1, \dots, \mathbf{x}_K] \in \mathbb{R}^{D \times K}$ where $d(\mathbf{x}_i, \mathbf{x}_k) < \tau, \forall k = 1, \dots, K$, $\exists \mathbf{w} \in \mathbb{R}^K$, s.t. $\mathbf{x}_i \approx \mathbf{X}_i \mathbf{w}_i$.

In this sense, given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, LLE aims at finding a low-dimensional embedding $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{L \times N}$ ($L < D$) that inherits the local linear self-representation relations.

Closed-form Solution for LLE

LLE can be decomposed into 3 steps

1. **Linear Reconstruction by Neighbors:** First, we compute the linear coefficients $\tilde{\mathbf{w}}$ by

$$\tilde{\mathbf{W}}^* = \arg \min_{\tilde{\mathbf{W}}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{X}_i \tilde{\mathbf{w}}_i\|_2^2 \quad \text{s.t.} \quad \tilde{\mathbf{W}} \mathbf{1}_K = \mathbf{1}_N$$

Here $\tilde{\mathbf{W}} = [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_N]^\top \in \mathbb{R}^{N \times K}$. The coefficient $\tilde{\mathbf{w}}_i = [\tilde{w}_{i1}, \dots, \tilde{w}_{iK}]^\top$ for each sample is constrained such that coefficients weighted on each neighbor sums up to 1. \mathbf{x}_i refers to the 'sample'

and \mathbf{X}_i refers to its 'neighbors'.

2. **Linear Embedding:** First we expand the old $\tilde{\mathbf{W}} = [\tilde{w}_{ij}] \in \mathbb{R}^{N \times K}$ to $\mathbf{W} = [w_{ij}] \in \mathbb{R}^{N \times N}$ by

$$w_{ij} = \begin{cases} \tilde{w}_{ij} & \text{if } \mathbf{x}_j \in \text{KNN}(\mathbf{x}_i) \\ 0 & \text{otherwise} \end{cases}$$

Compute the embedding $\mathbf{Z} \in \mathbb{R}^{L \times N}$ by

$$\mathbf{Z}^* = \arg \min_{\mathbf{Z}} \sum_{i=1}^N \left\| \mathbf{z}_i - \sum_{j=1}^N w_{ij} \mathbf{z}_j \right\|_2^2 \quad \text{s.t.} \quad \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i \mathbf{z}_i^\top = \mathbf{I}_L, \sum_{i=1}^N \mathbf{z}_i = \mathbf{0}$$

We constraint the embedding to ensure that $\text{Cov}(\mathbf{Z}) = \mathbf{I}_L$. The second constraint can be temporarily ignored since it can be achieved implicitly. We want to rewrite the object function in a more compact form.

$$\begin{aligned} \sum_{i=1}^N \left\| \mathbf{z}_i - \sum_{j=1}^N w_{ij} \mathbf{z}_j \right\|_2^2 &= \sum_{i=1}^N \left\| \mathbf{z}_i - \mathbf{Z} \mathbf{w}_i \right\|_2^2 = \sum_{i=1}^N \left\| \mathbf{Z} \mathbf{1}_i - \mathbf{Z} \mathbf{w}_i \right\|_2^2 = \left\| \mathbf{Z} - \mathbf{Z} \mathbf{W}^\top \right\|_F^2 \\ &= \text{tr} \left((\mathbf{Z} - \mathbf{Z} \mathbf{W}^\top)(\mathbf{Z} - \mathbf{Z} \mathbf{W}^\top)^\top \right) \\ &= \text{tr} \left(\mathbf{Z}(\mathbf{I} - \mathbf{W} - \mathbf{W}^\top + \mathbf{W}^\top \mathbf{W}) \mathbf{Z}^\top \right) \end{aligned}$$

where the **alignment matrix** $\Phi = \mathbf{I}_N - \mathbf{W} - \mathbf{W}^\top + \mathbf{W}^\top \mathbf{W}$.

3. Conduct EVD on $\Phi := \mathbf{U} \Lambda \mathbf{U}^\top$. After sorting the eigenvectors from smallest to largest eigenvalues, **we ignore the first eigenvector having zero eigenvalue** and take the L smallest eigenvectors of \mathbf{U} with non-zero eigenvalues as the embedding $(\mathbf{Z}^\top)^* \in \mathbb{R}^{N \times L}$.

First, for the **linear reconstruction by neighbors**, the coefficients \mathbf{W} can be computed as follows: Note that

$$\begin{aligned} \left\| \mathbf{x}_i - \mathbf{X}_i \mathbf{w}_i \right\|_2^2 &= \left\| \mathbf{x}_i (\mathbf{1}_K^\top \mathbf{w}_i) - \mathbf{X}_i \mathbf{w}_i \right\|_2^2 = \left\| (\mathbf{x}_i \mathbf{1}_K^\top - \mathbf{X}_i) \mathbf{w}_i \right\|_2^2 \\ &= \mathbf{w}_i^\top (\mathbf{x}_i \mathbf{1}_K^\top - \mathbf{X}_i)^\top (\mathbf{x}_i \mathbf{1}_K^\top - \mathbf{X}_i) \mathbf{w}_i \\ &\equiv \mathbf{w}_i^\top \mathbf{G}_i \mathbf{w}_i \end{aligned}$$

where we denote $\mathbf{G}_i = (\mathbf{x}_i \mathbf{1}_K^\top - \mathbf{X}_i)^\top (\mathbf{x}_i \mathbf{1}_K^\top - \mathbf{X}_i) \in \mathbb{R}^{K \times K}$. The optimization problem is

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \sum_{i=1}^N \mathbf{w}_i^\top \mathbf{G}_i \mathbf{w}_i \quad \text{s.t.} \quad \mathbf{W} \mathbf{1}_K = \mathbf{1}_N$$

The Lagrangian for this is

$$\begin{aligned} \mathcal{L}(\mathbf{W}, \Lambda) &= \sum_{i=1}^N \mathbf{w}_i^\top \mathbf{G}_i \mathbf{w}_i - \sum_{i=1}^N \lambda_i (\mathbf{1}_K^\top \mathbf{w}_i - 1) \\ \Rightarrow \begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_i} = 2 \mathbf{G}_i \mathbf{w}_i - \lambda_i \mathbf{1}_K = \mathbf{0} \\ \frac{\partial \mathcal{L}}{\partial \lambda_i} = \mathbf{1}_K^\top \mathbf{w}_i - 1 = 0 \end{cases} &\Rightarrow \begin{cases} \mathbf{w}_i = \frac{\lambda_i}{2} \mathbf{G}_i^{-1} \mathbf{1}_K \\ \mathbf{1}_K^\top \frac{\lambda_i}{2} \mathbf{G}_i^{-1} \mathbf{1}_K = 1 \end{cases} \\ \Rightarrow \begin{cases} \mathbf{w}_i = \frac{\lambda_i}{2} \mathbf{G}_i^{-1} \mathbf{1}_K \\ \lambda_i = \frac{2}{\mathbf{1}_K^\top \mathbf{G}_i^{-1} \mathbf{1}_K} \end{cases} &\Rightarrow \mathbf{w}_i = \frac{\mathbf{G}_i^{-1} \mathbf{1}_K}{\mathbf{1}_K^\top \mathbf{G}_i^{-1} \mathbf{1}_K} \end{aligned}$$

Second, for the derivation of **linear embedding**, our optimization problem is essentially

$$\min \text{tr}(\mathbf{Z}\Phi\mathbf{Z}^\top) \quad \text{s.t.} \quad \frac{1}{N}\mathbf{Z}\mathbf{Z}^\top = \mathbf{I}_L$$

and therefore the Lagrangian for this is (**Important: Under optimal** $\mathbf{\Lambda} \in \mathbb{R}^{L \times L}$)

$$\begin{aligned} \mathcal{L}(\mathbf{Z}, \mathbf{\Lambda}) &= \text{tr}(\mathbf{Z}\Phi\mathbf{Z}^\top) - \text{tr}(\mathbf{\Lambda}^\top (\frac{1}{N}\mathbf{Z}\mathbf{Z}^\top - \mathbf{I}_L)) \\ \Rightarrow \frac{\partial \mathcal{L}}{\partial \mathbf{Z}} &= 2\mathbf{Z}\Phi - \frac{2}{N}\mathbf{\Lambda}\mathbf{Z} = 0 \Rightarrow \Phi\mathbf{Z}^\top = \mathbf{Z}^\top (\frac{1}{N}\mathbf{\Lambda}) \end{aligned}$$

Moreover, recall that our goal is to minimize

$$\text{tr}(\mathbf{Z}\Phi\mathbf{Z}^\top) = \text{tr}(\mathbf{Z}\mathbf{Z}^\top \frac{1}{N}\mathbf{\Lambda}) = \text{tr}(\frac{1}{N}\mathbf{\Lambda}) = \frac{1}{N} \sum_{i=1}^N \lambda_i$$

and EVD of $\Phi := \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$. This means that **under optimal**, we **should pick L eigenvectors** from the eigenvectors of Φ to compose the embedding $(\mathbf{Z}^\top)^* \in \mathbb{R}^{N \times L}$. After sorting the eigenvectors from smallest to largest eigenvalues, **we ignore the first eigenvector having zero eigenvalue** and take the L smallest eigenvectors of \mathbf{U} with non-zero eigenvalues of Φ as the embedding $(\mathbf{Z}^\top)^*$.

1.4 Laplacian Eigenmap

Laplacian Eigenmap

Given a set of data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, we construct the similarity matrix $\mathbf{A} = [a(\mathbf{x}_i, \mathbf{x}_j)] \in \mathbb{R}^{N \times N}$. A reasonable criterion to get the low-dimensional embedding $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{L \times N}$ is to minimize the following objective function

$$\min_{\mathbf{Z}} \sum_{m,n=1}^N \|\mathbf{z}_m - \mathbf{z}_n\|_2^2 a(\mathbf{x}_m, \mathbf{x}_n)$$

because when distance $\|\mathbf{z}_m - \mathbf{z}_n\|_2^2$ is small, the similarity $a(\mathbf{x}_m, \mathbf{x}_n)$ should be large.

Closed-form Solution of Laplacian Eigenmap

$$\begin{aligned} \mathbf{Z} &= \arg \min_{\mathbf{Z}} \sum_{m,n=1}^N \|\mathbf{z}_m - \mathbf{z}_n\|_2^2 a(\mathbf{x}_m, \mathbf{x}_n) = \arg \min_{\mathbf{Z}} \sum_{m,n=1}^N (\mathbf{z}_m^\top \mathbf{z}_m - 2\mathbf{z}_m^\top \mathbf{z}_n + \mathbf{z}_n^\top \mathbf{z}_n) a_{mn} \\ &= \arg \min_{\mathbf{Z}} \sum_{m=1}^N \mathbf{z}_m^\top \mathbf{z}_m \left(\sum_{n=1}^N a_{mn} \right) + \sum_{n=1}^N \mathbf{z}_n^\top \mathbf{z}_n \left(\sum_{m=1}^N a_{mn} \right) - 2 \sum_{m,n=1}^N \mathbf{z}_m^\top \mathbf{z}_n a_{mn} \\ &= \arg \min_{\mathbf{Z}} 2\text{tr}(\mathbf{Z}^\top \text{diag}(\mathbf{A}\mathbf{1}_N)\mathbf{Z}) - 2\text{tr}(\mathbf{Z}^\top \mathbf{A}\mathbf{Z}) \\ &= \arg \min_{\mathbf{Z}} 2\text{tr}(\mathbf{Z}^\top (\text{diag}(\mathbf{A}\mathbf{1}_N) - \mathbf{A})\mathbf{Z}) \\ &= \arg \min_{\mathbf{Z}} \text{tr}(\mathbf{Z}^\top \mathbf{L}\mathbf{Z}) \quad \text{where } \mathbf{L} = \text{diag}(\mathbf{A}\mathbf{1}_N) - \mathbf{A} \end{aligned}$$

In practice, the Laplacian matrix \mathbf{L} is usually normalized by the degree matrix $\mathbf{D} = \text{diag}(\mathbf{A}\mathbf{1}_N)$:

$$\hat{\mathbf{L}}_{\text{sym}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}} = \mathbf{D}^{-\frac{1}{2}} (\mathbf{D} - \mathbf{A}) \mathbf{D}^{-\frac{1}{2}} = \mathbf{I}_N - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} = \mathbf{I}_N - \hat{\mathbf{A}}$$

By performing EVD on $\hat{\mathbf{L}}_{\text{sym}} := \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, we can get the embedding $\mathbf{Z}^* = \mathbf{U}_L \in \mathbb{R}^{N \times L}$.

In construction of similarity matrix \mathbf{A} , we can apply the Gram matrix of kernel function such as the RBF kernel:

$$a(\mathbf{x}_i, \mathbf{x}_j) := K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2/h)$$

2 Kernel Methods

Kernel PCA

Suppose our data $\mathbf{X} \in \mathbb{R}^{N \times D}$ is non-linearly separable. We can first map the data into a higher-dimensional space $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]^\top \in \mathbb{R}^{N \times \dim(F)}$ and then perform EVD on the Gram matrix $\mathbf{K} = \Phi(\mathbf{X})\Phi(\mathbf{X})^\top$.

$$\mathbf{K} = \mathbf{V}\mathbf{\Delta}\mathbf{V}^\top$$

The PCA corresponds to the top- L eigenvectors of \mathbf{K} : $\mathbf{Z}^* = \mathbf{V}_L \mathbf{\Delta}_L^{\frac{1}{2}} \in \mathbb{R}^{N \times L}$.

Revisiting MDS and ISOMAP, we can consider them as special cases of Kernel PCA.

- For MDS, $\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C} = \mathbf{C}\mathbf{X}\mathbf{X}^\top\mathbf{C} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$ (Linear Kernel)
- For ISOMAP, $\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D}_{geo} \odot \mathbf{D}_{geo})\mathbf{C}$ (Mercer Kernel)

2.1 t-Distributed Stochastic Neighbor Embedding (t-SNE)

t-SNE

Given a dataset $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, first we define a Probability p_{ij} that is proportional to the similarity between \mathbf{x}_i and \mathbf{x}_j :

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N}, \quad p_{ii} = 0$$

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|_2^2/2\sigma_i^2)}$$

t-SNE aims to learn $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{K \times N}$ (usually $K = 2, 3$ for visualization purposes) that minimizes the KL divergence between p_{ij} and q_{ij}

$$\min_{\mathbf{Z}} \text{KL}(P||Q) = \min_{\mathbf{Z}} \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

where q_{ij} is the similarity between \mathbf{z}_i and \mathbf{z}_j :

$$q_{ij} = \frac{(1 + \|\mathbf{z}_i - \mathbf{z}_j\|_2^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{z}_k - \mathbf{z}_l\|_2^2)^{-1}}, \quad q_{ii} = 0$$

where $\{q_{ij}\}$ is the Student-t distribution with $\text{df}=1$. Optimization of KL divergence is done with SGD.

3 Autoencoding

First, let's revisit PCA from a viewpoint of **autoencoding**. Recall that PCA is the least-square data denoising under i.i.d. Gaussian noise,

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \Omega} \|\mathbf{X}_{\text{noisy}} - \mathbf{X}\|_F^2 = \mathbf{U}_L \mathbf{\Sigma}_L \mathbf{V}_L^\top, \text{ where } \mathbf{X}_{\text{noisy}} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$$

referring to the construction of principal components and the corresponding reconstruction. This can be viewed as a special case of autoencoding where the encoder and decoder are linear transformations.

$$\text{Encoder: } \mathbf{Z} = \mathbf{X}_{\text{noisy}} \mathbf{V}_L^\top$$

$$\text{Decoder: } \mathbf{X}^* = \mathbf{X}_{\text{noisy}} \mathbf{V}_L^\top \mathbf{V}_L$$

Here \mathbf{V}_L^\top and \mathbf{V}_L are the encoder and decoder respectively.

Autoencoders

In general, a typical autocoder consists of

$$\text{Encoder: } f : \mathcal{X} \rightarrow \mathcal{Z}$$

$$\text{Decoder: } g : \mathcal{Z} \rightarrow \mathcal{X}$$

Given a set of data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{D \times N}$, the autoencoder aims to learn the encoder and decoder that minimize the reconstruction error

$$\min_{f,g} \sum_{i=1}^N \text{loss}(\mathbf{x}_i - g(f(\mathbf{x}_i))) + \text{regularization}(q_{\mathbf{Z}|\mathbf{X}}, p_{\mathbf{Z}})$$

where $q_{\mathbf{Z}|\mathbf{X}}$ is the posterior distribution of latent space \mathbf{Z} given dataset \mathbf{X} and $p_{\mathbf{Z}}$ is the prior distribution of \mathbf{Z} .

References

- [Multidimensional Scaling, Sammon Mapping, and Isomap: Tutorial and Survey](#)
- [Locally Linear Embedding and its Variants: Tutorial and Survey](#)

Appendix

Classic MDS

The specific process of deriving $\mathbf{K} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$ is as follows: Note that

$$\begin{aligned}\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N &= \begin{bmatrix} x_{11}^2 & \cdots & x_{1D}^2 \\ \vdots & \ddots & \vdots \\ x_{N1}^2 & \cdots & x_{ND}^2 \end{bmatrix}_{N \times D} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{D \times 1} \begin{bmatrix} 1 & \cdots & 1 \end{bmatrix}_{1 \times N} = \begin{bmatrix} \mathbf{x}_1^\top \mathbf{x}_1 & \cdots & \mathbf{x}_1^\top \mathbf{x}_1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_N^\top \mathbf{x}_N & \cdots & \mathbf{x}_N^\top \mathbf{x}_N \end{bmatrix}_{N \times N} \\ &= \begin{bmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{x}_1 \rangle \\ \vdots & \ddots & \vdots \\ \langle \mathbf{x}_N, \mathbf{x}_N \rangle & \cdots & \langle \mathbf{x}_N, \mathbf{x}_N \rangle \end{bmatrix}_{N \times N}\end{aligned}$$

and

$$\begin{aligned}d_{ij} &= \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_i - 2\mathbf{x}_i^\top \mathbf{x}_j + \mathbf{x}_j^\top \mathbf{x}_j \\ &= \langle \mathbf{x}_i, \mathbf{x}_i \rangle + \langle \mathbf{x}_j, \mathbf{x}_j \rangle - 2\langle \mathbf{x}_i, \mathbf{x}_j \rangle\end{aligned}$$

We can decompose

$$\mathbf{D} \odot \mathbf{D} = \begin{bmatrix} d_{11}^2 & \cdots & d_{1N}^2 \\ \vdots & \ddots & \vdots \\ d_{N1}^2 & \cdots & d_{NN}^2 \end{bmatrix} = (\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N) + (\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N)^\top - 2\mathbf{X}\mathbf{X}^\top$$

\mathbf{C} is essentially a **centering matrix** since

$$\begin{aligned}\mathbf{C}\mathbf{X} &= (\mathbf{I}_N - \frac{1}{N}\mathbf{1}_{N \times N})\mathbf{X} = \mathbf{X} - \frac{1}{N}\mathbf{1}_N \mathbf{1}_N^\top \mathbf{X} \\ &= \begin{bmatrix} x_{11} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{ND} \end{bmatrix} - \begin{bmatrix} \frac{x_{11} + \cdots + x_{N1}}{N} & \cdots & \frac{x_{1D} + \cdots + x_{ND}}{N} \\ \vdots & \ddots & \vdots \\ \frac{x_{11} + \cdots + x_{N1}}{N} & \cdots & \frac{x_{1D} + \cdots + x_{ND}}{N} \end{bmatrix} \\ &= \begin{bmatrix} x_{11} - \bar{x}_1 & \cdots & x_{1D} - \bar{x}_D \\ \vdots & \ddots & \vdots \\ x_{N1} - \bar{x}_1 & \cdots & x_{ND} - \bar{x}_D \end{bmatrix}\end{aligned}$$

Therefore, $\tilde{\mathbf{X}} = \mathbf{C}\mathbf{X}$ is the zero-meaned data matrix of \mathbf{X} . One can verify that after **double centralizing**,

$$\mathbf{C}((\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N + (\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N)^\top))\mathbf{C} = \mathbf{0}_{N \times N}$$

Therefore, the inner product data \mathbf{K} is essentially

$$\begin{aligned}\mathbf{K} &= -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C} = -\frac{1}{2}\mathbf{C}(\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N + (\mathbf{X} \odot \mathbf{X} \mathbf{1}_D \mathbf{1}_N)^\top - 2\mathbf{X}\mathbf{X}^\top)\mathbf{C} \\ &= \mathbf{C}\mathbf{X}\mathbf{X}^\top\mathbf{C} = \mathbf{C}\mathbf{X}(\mathbf{C}\mathbf{X})^\top = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\end{aligned}$$