

Some proof on Lecture 6-8

Lin Ning linning51400@ruc.edu.cn

2022/4/23

Reference

Source 1 Mathematics for Machine Learning by Marc Peter Deisenroth A Aldo Faisal Cheng Soon Ong, pp. 131-132, 326-331

Source 2 机器学习，周志华， pp. 228-233

Source 3 www.stat.yale.edu/~lc436/papers/JCGS-mds.pdf

Transformation

LLE

In step 1 we choose k-nearest point of \mathbf{x}_i and find their best linear representation in order to reduce the related point and get sparse weight matrix, the optimal target is

$$\begin{aligned} \mathbf{W}^* &:= \arg \min_{\mathbf{W}} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{j \in \mathcal{N}_K(\mathbf{x}_i)} w_{ij} \mathbf{x}_j \right\|_2^2, \\ s. t. \quad &\sum_{j=1}^K w_{ij} = \sum_{j=1}^N w_{ij} = 1 \Rightarrow \mathbf{W} \mathbf{1}_N = \mathbf{1}_N \end{aligned}$$

Define K-NN matrix \mathbf{N}_i and their weight matrix of \mathbf{x}_i :

$$\mathbf{N}_i := \{\mathbf{x}_j\}_{j \in \mathcal{N}_K(\mathbf{x}_i)} \in \mathbb{R}^{D \times K}$$

and the optimal solution is [1]

$$\boldsymbol{\omega}_i = \{w_{ij}\}_{j \in \mathcal{N}_K(\mathbf{x}_i)}^T \in \mathbb{R}^{K \times 1}$$

The optimal problem is

$$\boldsymbol{\omega}_i := \arg \min_{\boldsymbol{\omega}} \|\mathbf{x}_i - \mathbf{N}_i \boldsymbol{\omega}\|_2^2, \text{ s.t. } \mathbf{1}_K^T \boldsymbol{\omega} = 1, i = 1, 2, \dots, N$$

Assume that \mathbf{x}_i and its K nearest vectors are linearly independent, let

$$\mathbf{C} := (\mathbf{N}_i - \mathbf{x}_i \mathbf{1}_K^T)^T (\mathbf{N}_i - \mathbf{x}_i \mathbf{1}_K^T)$$

and

$$\boldsymbol{\omega}^* = \frac{\mathbf{C}^{-1} \mathbf{1}_K}{\mathbf{1}_K^T \mathbf{C}^{-1} \mathbf{1}_K} = \text{norm}(\mathbf{C}^{-1} \mathbf{1}_K)$$

where norm is the function of normalization through sum all the elements of the vector

$$\text{norm}(\mathbf{x}) := \frac{1}{\sum_i x_i} \mathbf{x}$$

Using the parameters in $\boldsymbol{\omega}_i$ to fill the element w_{ij} of \mathbf{W} , step 1 end.

In step 2 we use the low-dimension vector \mathbf{z}_i to get the embedding of weight matrix. In order to avoid trivial solution, the optimal target is

$$\begin{aligned} \mathbf{Z}^* &:= \arg \min_{\mathbf{Z}} \sum_{i=1}^N \left\| \mathbf{z}_i - \sum_{j=1}^N w_{ij} \mathbf{z}_j \right\|_2^2 = \arg \min_{\mathbf{Z}} \left\| \mathbf{Z} - \mathbf{Z} \mathbf{W}^T \right\|_F^2 \\ &= \arg \min_{\mathbf{Z}} \left\| \mathbf{Z}^T - \mathbf{W} \mathbf{Z}^T \right\|_F^2, \text{ s.t. } \mathbf{Z} \mathbf{Z}^T = \mathbf{I}_D, \mathbf{Z} \mathbf{1}_N = \mathbf{0}_N \end{aligned}$$

The shifting of \mathbf{z}_i won't affect the optimal function [2], i.e. won't help to avoid trivial solution, so the zero-mean constraint always be discarded.

Let

$$\boldsymbol{\Phi} := (\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})$$

and

$$\mathbf{Z}^* := \arg \min_{\mathbf{Z}} \text{tr}(\mathbf{Z} \boldsymbol{\Phi} \mathbf{Z}^T), \text{ s.t. } \mathbf{Z} \mathbf{Z}^T = \mathbf{I}_D$$

PCA

Assuming that data matrix $\mathbf{X} \in \mathbb{R}^{D \times N}$.

The perspective of projection try to minimize the reconstruction error of in dimension reduction by setting the optimal target as follows

$$\begin{aligned}\mathbf{X}^* &:= \arg \min_{\tilde{\mathbf{X}}} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}) = \arg \min_{\tilde{\mathbf{X}}} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2^2 \\ &= \arg \min_{\tilde{\mathbf{X}}} \left\| \mathbf{X} - \tilde{\mathbf{X}} \right\|_F^2, \text{ s.t. } \text{rank}(\tilde{\mathbf{X}}) \leq L\end{aligned}$$

or

$$\mathbf{X}^* = \arg \min_{\tilde{\mathbf{X}} \in \Omega} \left\| \mathbf{X} - \tilde{\mathbf{X}} \right\|_F^2, \text{ s.t. } \dim(\Omega) \leq L$$

Select an unit orthogonal basis for L dimensional principle space

$$\tilde{\mathbf{x}}_n := \sum_{i=1}^L z_{in} \mathbf{b}_i = \mathbf{B} \mathbf{z}_n$$

and

$$z_{in} = \langle \tilde{\mathbf{x}}_n, \mathbf{b}_i \rangle = \tilde{\mathbf{x}}_n^T \mathbf{b}_i = \mathbf{b}_i^T \tilde{\mathbf{x}}_n$$

and

$$\mathbf{z}_n = \mathbf{B}^T \tilde{\mathbf{x}}_n$$

Therefore

$$\tilde{\mathbf{X}} = \mathbf{B} \mathbf{Z}$$

and

$$\mathbf{Z} = \mathbf{B}^T \tilde{\mathbf{X}}$$

PCA actually define an orthogonal projection from the primal space to principal space [3], i.e.

$$\tilde{\mathbf{X}} = \mathbf{B} \mathbf{B}^T \mathbf{X}$$

and

$$\mathbf{Z} = \mathbf{B}^T \mathbf{B} \mathbf{B}^T \mathbf{X} = \mathbf{B}^T \mathbf{X}$$

The optimal target can be transformed

$$\begin{aligned} \mathbf{X}^* &= \arg \min_{\mathbf{B}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{B} \mathbf{z}_n\|_2^2 = \arg \min_{\mathbf{B}} \|\mathbf{X} - \mathbf{B} \mathbf{Z}\|_F^2 \\ &= \arg \min_{\mathbf{B}} \text{tr}((\mathbf{X} - \mathbf{B} \mathbf{Z})^T (\mathbf{X} - \mathbf{B} \mathbf{Z})) \\ &= \arg \min_{\mathbf{B}} \text{tr}(\mathbf{Z}^T \mathbf{Z}) - 2\text{tr}(\mathbf{Z}^T \mathbf{B}^T \mathbf{X}) = \arg \min_{\mathbf{B}} -\text{tr}(\mathbf{Z}^T \mathbf{Z}) \\ &= \arg \max_{\mathbf{B}} \text{tr}(\mathbf{Z}^T \mathbf{Z}) = \arg \max_{\mathbf{B}} \text{tr}(\mathbf{Z} \mathbf{Z}^T) \\ &= \arg \max_{\mathbf{B}} \text{tr}(\mathbf{B}^T \mathbf{X} \mathbf{X}^T \mathbf{B}), \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_L \end{aligned}$$

Let

$$\Phi := \sum_{j=L+1}^D \mathbf{x}_j \mathbf{x}_j^T = \mathbf{X} \mathbf{X}^T \in \mathbf{R}^{N \times N}$$

The optimal target is

$$\mathbf{B}^* = \arg \max_{\mathbf{B}} \text{tr}(\mathbf{B}^T \Phi \mathbf{B}), \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_L$$

or

$$\mathbf{B}^* = \arg \min_{\mathbf{B}} -\text{tr}(\mathbf{B}^T \Phi \mathbf{B}), \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_L$$

therefore

$$\mathbf{X}^* = \mathbf{B}^* \mathbf{B}^{*T} \mathbf{X}$$

To compress the data, we should use reversible linear maps to transform the data from the primal space to principal space. Since the change of base will preserve the information of data, using the orthogonal base of principle space, we represent data as

$$\mathbf{X}_p := \mathbf{Z}^* = \mathbf{B}^{*T} \mathbf{X}$$

The \mathbf{B} can be thought of encoder and the \mathbf{B}^T can be thought of decoder.

Kernel PCA

for kernel PCA, we apply PCA in kernel space

$$\mathbf{X}^* = \arg \min_{\tilde{\mathbf{X}}} \left\| \phi(\mathbf{X}) - \tilde{\mathbf{X}} \right\|_F^2, \text{ s.t. } \text{rank}(\tilde{\mathbf{X}}) \leq L$$

make the substitution and we derive

$$\Phi := \sum_{j=1}^N \phi(\mathbf{x}_j) \phi(\mathbf{x}_j)^T = \phi(\mathbf{X}) \phi(\mathbf{X})^T \in \mathbf{R}^{N \times N}$$

and

$$\mathbf{B}^* = \arg \max_{\mathbf{B}} \text{tr}(\mathbf{B}^T \Phi \mathbf{B}), \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_L$$

Which we pay more attention to is the representation in low-dimensional space

$$\mathbf{X}_p := \mathbf{Z}^* = \mathbf{B}^{*T} \mathbf{X}$$

EigenMap

Define the similarity function and matrix

$$a(\mathbf{x}_n, \mathbf{x}_m) = a(\mathbf{x}_m, \mathbf{x}_n) \in [0, 1]$$

and

$$\mathbf{A} = \{a(\mathbf{x}_n, \mathbf{x}_m)\}_{n,m=1}^D$$

For eigenmap, the optimal target is

$$\begin{aligned}
\mathbf{Z}^* &:= \arg \min_{\mathbf{Z}} \sum_{n,m=1}^D \|\mathbf{z}_n - \mathbf{z}_m\|_2^2 a_{nm} \\
&= \arg \min_{\mathbf{Z}} \sum_{n,m=1}^D (z_n^T z_n - 2z_n^T z_m + z_m^T z_m) a_{nm} \\
&= \arg \min_{\mathbf{Z}} \left(\sum_{n=1}^D z_n^T z_n \sum_{m=1}^D a_{nm} + \sum_{m=1}^D z_m^T z_m \sum_{n=1}^D a_{nm} - 2 \sum_{n,m=1}^D z_n^T z_m a_{nm} \right) \\
&= \arg \min_{\mathbf{Z}} (2\text{tr}(\mathbf{Z}^T \text{diag}(\mathbf{A}\mathbf{1}_N)\mathbf{Z}) - 2\text{tr}(\mathbf{Z}^T \mathbf{A}\mathbf{Z})), \text{ s.t. } \mathbf{Z}^T \mathbf{Z} = \mathbf{I}_N
\end{aligned}$$

Let

$$\Phi = \text{diag}(\mathbf{A}\mathbf{1}_D) - \mathbf{A}$$

and

$$\mathbf{Z}^* = \arg \min_{\mathbf{Z}} \text{tr}(\mathbf{Z}^T \Phi \mathbf{Z}), \text{ s.t. } \mathbf{Z}^T \mathbf{Z} = \mathbf{I}_L$$

ISOMap

Define the distance function and matrix

$$d(\mathbf{x}_n, \mathbf{x}_m) = d(\mathbf{x}_m, \mathbf{x}_n) \geq 0$$

and

$$\mathbf{D} = \{d(\mathbf{x}_n, \mathbf{x}_m)\}_{n,m=1}^N$$

In order to find a low-dimensional representation in latent space whose Euclidean distance is close to the defined distance and get the best representation, the optimal target of Classic MDS is [4]

$$\begin{aligned}
\mathbf{Z}^* &= \arg \min_{\mathbf{Z}} \text{Strain}_d(\mathbf{Z}) = \arg \min_{\mathbf{Z}} \sum_{n,m=1}^N (k_{nm} - z_n^T z_m)^2 \\
&= \arg \min_{\mathbf{Z}} \|\mathbf{K} - \mathbf{Z}^T \mathbf{Z}\|_F^2, \text{ s.t. } \text{rank}(\mathbf{Z}) \leq L
\end{aligned}$$

where

$$\mathbf{K} = -\frac{1}{2}\mathbf{C}(\mathbf{D} \odot \mathbf{D})\mathbf{C}$$

and

$$\mathbf{C} = \mathbf{I}_N - \frac{1}{N}\mathbf{1}_{N \times N}$$

In ISOMAP, we always make assumption that \mathbf{K} is positive semi-definite or force \mathbf{K} to be positive semi-definite using the similar idea of ridge regression.

Define Gram matrix of \mathbf{Z}

$$\mathbf{G} := \{\mathbf{z}_n^T \mathbf{z}_m\}_{n,m=1}^N = \mathbf{Z}^T \mathbf{Z}, \text{ rank}(\mathbf{G}) = \text{rank}(\mathbf{Z})$$

Do eigenvalue decomposition to postive semi-definite positive \mathbf{G} , we get

$$\mathbf{G} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$$

and

$$\mathbf{Z} = \mathbf{\Lambda}^{1/2} \mathbf{P}^T$$

The optimal problem can be converted to PCA

$$\mathbf{G}^* = \arg \min_{\mathbf{G}} \|\mathbf{K} - \mathbf{G}\|_F^2, \text{ s.t. } \text{rank}(\mathbf{G}) \leq L, \mathbf{G}^T = \mathbf{G}, \mathbf{x}^T \mathbf{G} \mathbf{x} \geq 0$$

According to PCA, the $\mathbf{\Phi}$ is define as

$$\mathbf{\Phi} := \sum_{j=1}^N \mathbf{k}_j \mathbf{k}_j^T = \mathbf{K} \mathbf{K}^T \in \mathbf{R}^{N \times N}$$

and

$$\mathbf{B}^* = \arg \max_{\mathbf{B}} \text{tr}(\mathbf{B}^T \mathbf{\Phi} \mathbf{B}), \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_N$$

and

$$\mathbf{G}^* = \mathbf{B}_L^* \mathbf{B}_L^{*T} \mathbf{K}$$

The restriction of \mathbf{G} can be satisfied. [5]

Solution

Major Problem

It's easy to conclude that all the Φ we defined is symmetric and positive semi-definite [6].

Assume \mathbf{Z} is a $D \times L$ matrix and suppose $D > L$.

The next step is how to solve the optimal problem as

$$\mathbf{Z}^* = \arg \min_{\mathbf{Z}} \text{tr}(\mathbf{Z}^T \Phi \mathbf{Z}), \text{ s.t. } \mathbf{Z}^T \mathbf{Z} = \mathbf{I}_D$$

Apply eigenvalue decomposition to Φ to simplify the problem

$$\Phi = \mathbf{V} \Lambda \mathbf{V}^T, \mathbf{V} \mathbf{V}^T = \mathbf{I}_D$$

Sort the eigenvalue to the satisfies following equation

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$$

And rearrange the corresponding eigenvector

$$\Phi \mathbf{v}_i = \lambda_i \mathbf{v}_i, i = 1, 2, \dots, D$$

Let

$$\mathbf{Y} = \mathbf{V}^T \mathbf{Z}, \mathbf{Y}^T \mathbf{Y} = \mathbf{Z}^T \mathbf{V} \mathbf{V}^T \mathbf{Z} = \mathbf{I}_D$$

and

$$\mathbf{Z} = \mathbf{V} \mathbf{Y}$$

Using the the bijection, we can simplify the optimal problem

$$\mathbf{Y}^* = \arg \min_{\mathbf{Y}} \text{tr}(\mathbf{Y}^T \Lambda \mathbf{Y}) = \arg \min_{\mathbf{Y}} \sum_{n=1}^L \mathbf{y}_n^T \Lambda \mathbf{y}_n, \text{ s.t. } \mathbf{Y}^T \mathbf{Y} = \mathbf{I}_D$$

and

$$\mathbf{Z}^* = \mathbf{V} \mathbf{Y}^*$$

The optimal problem may difficult to directly solve for the complex constrains for us, so we should make full use of its excellent algebra property. Let

$$\mathbf{Y}^T = (\gamma_1, \gamma_2, \dots, \gamma_D), \gamma'_i \in \mathbb{R}^{L \times 1}$$

and

$$\begin{aligned} \mathbf{Y}^* &= \arg \min_{\mathbf{Y}} \text{tr}(\mathbf{Y}^T \mathbf{\Lambda} \mathbf{Y}) \\ &= \arg \min_{\mathbf{Y}} \text{tr} \left(\sum_{n=1}^D \lambda_n \gamma_n \gamma_n^T \right) \\ &= \arg \min_{\mathbf{Y}} \sum_{n=1}^D \lambda_n \text{tr}(\gamma_n \gamma_n^T) \end{aligned}$$

Let

$$a_n = \text{tr}(\gamma_n \gamma_n^T)$$

After analyze the constraints on a_n [7], we can prove that [8]

$$\sum_{n=1}^D \lambda_n a_n \geq \sum_{n=D-L+1}^D \lambda_n, \text{ s.t. } \sum_{n=1}^D a_n = L, 0 \leq a_i \leq 1, i = 1, 2, \dots, D$$

We finally conclude that

$$\min_{\mathbf{Z}} \text{tr}(\mathbf{Z}^T \mathbf{\Phi} \mathbf{Z}) = \min_{\mathbf{Y}} \text{tr}(\mathbf{Y}^T \mathbf{\Lambda} \mathbf{Y}) = \min_{\mathbf{Y}} \sum_{n=1}^D \lambda_n \text{tr}(\gamma_n \gamma_n^T) \geq \sum_{n=D-L+1}^D \lambda_n$$

Select eigenvectors corresponding to L th smallest eigenvalues to construct \mathbf{Z}

$$\mathbf{Z} = (\mathbf{v}_{D-L+1}, \mathbf{v}_{D-L+2}, \dots, \mathbf{v}_D)$$

and

$$\text{tr}(\mathbf{Z}^T \mathbf{\Phi} \mathbf{Z}) = \sum_{n=1}^L \mathbf{z}_n^T \mathbf{\Phi} \mathbf{z}_n = \sum_{n=1}^L \lambda_{D-L+n} \mathbf{z}_n^T \mathbf{z}_n = \sum_{n=D-L+1}^D \lambda_n$$

So the answer is

$$\mathbf{Z}^* = (\mathbf{v}_{D-L+1}, \mathbf{v}_{D-L+2}, \dots, \mathbf{v}_D) = \mathbf{V}_{-L}$$

and

$$\min_{\mathbf{Z}} \text{tr}(\mathbf{Z}^T \mathbf{\Phi} \mathbf{Z}) = \sum_{n=D-L+1}^D \lambda_n$$

If we change the minimizing problem into maximizing problem, through similarly derivation, the solution is

$$\mathbf{Z}^* = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_L) = \mathbf{V}_L$$

and

$$\max_{\mathbf{Z}} \text{tr}(\mathbf{Z}^T \mathbf{\Phi} \mathbf{Z}) = \sum_{n=1}^L \lambda_n$$

The solution is interesting for PCA, kernel PCA, and ISOMap when combined with SVD.

PCA

Apply SVD to \mathbf{X}

$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

The optimal solution and representation is

$$\mathbf{X}^* = \mathbf{U}_L \mathbf{U}_L^T \mathbf{X} = \mathbf{U}_L \mathbf{\Sigma}_L \mathbf{V}_L^T$$

and

$$\mathbf{X}_p = \mathbf{U}_L^T \mathbf{X} = \mathbf{\Sigma}_L \mathbf{V}_L^T$$

We can also get the optimal reconstruct error

$$\min_{\tilde{\mathbf{X}}} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}) = \sum_{n=1}^L \sigma_n$$

It can be proved that the solution is also the optimal one in the following optimal problem [8]

$$\mathbf{X}^* = \arg \min_{\tilde{\mathbf{X}} \in \Omega} \left\| \mathbf{X} - \tilde{\mathbf{X}} \right\|_2, \text{ s. t. } \dim(\Omega) \leq L$$

where the matrix norm $\| \cdot \|_2$ is spectral norm defined as

$$\|A\|_2 := \max_{\mathbf{x}} \frac{\|A\mathbf{x}\|_2}{\|\mathbf{x}\|_2}$$

Kernel PCA

Apply SVD to $\phi(\mathbf{X})$ and we get

$$\phi(\mathbf{X}) = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

and the representation

$$\mathbf{X}_p = \mathbf{U}_L^T \phi(\mathbf{X})$$

In fact, the specific form of ϕ is always unknown so we define kernel function and its matrix

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

and

$$\mathbf{K} = \{\kappa(\mathbf{x}_n, \mathbf{x}_m)\}_{n,m=1}^N = \phi(\mathbf{X})^T \phi(\mathbf{X})$$

It will be shown in the derivation of the optimal problem

$$\mathbf{B}^* = \arg \max_{\mathbf{B}} \text{tr}(\mathbf{B}^T \mathbf{\Phi} \mathbf{B}), \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_L$$

We can derive the following result in a winding way [8], in fact, we can take full use of SVD. Consider the relationship between SVD and eigenvalue decomposition, \mathbf{K} , we get

$$\mathbf{K} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$$

So we can solve the problem using matrix \mathbf{K}

$$\mathbf{X}_p = \mathbf{U}_L^T \phi(\mathbf{X}) = \mathbf{\Sigma}_L \mathbf{V}_L^T$$

although there is no explicit solution to \mathbf{X}^*

ISOMap

Apply eigenvalue decomposition to \mathbf{K}

$$\mathbf{K} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$$

and the solution is

$$\mathbf{G}^* = \mathbf{V}_L \mathbf{V}_L^T \mathbf{K} = \mathbf{V}_L \mathbf{\Lambda}_L \mathbf{V}_L^T$$

Consider the relationship between SVD and eigenvalue decomposition, we get

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

Appendix

[1]

According to the constraints

$$\begin{aligned}\boldsymbol{\omega}_i &= \arg \min_{\boldsymbol{\omega}} \left\| \mathbf{x}_i \mathbf{1}_K^T \boldsymbol{\omega} - \mathbf{N}_i \boldsymbol{\omega} \right\|_2^2 \\ &= \arg \min_{\boldsymbol{\omega}} \left\| (\mathbf{x}_i \mathbf{1}_K^T - \mathbf{N}_i) \boldsymbol{\omega} \right\|_2^2, \text{ s.t. } \mathbf{1}_K^T \boldsymbol{\omega} = 1\end{aligned}$$

Let

$$\mathbf{Y} := \mathbf{x}_i \mathbf{1}_K^T - \mathbf{N}_i$$

and

$$\boldsymbol{\omega}_i = \arg \min_{\boldsymbol{\omega}} \left\| \mathbf{Y} \boldsymbol{\omega} \right\|_2^2, \text{ s.t. } \mathbf{1}_K^T \boldsymbol{\omega} = 1$$

Construct Lagrangian

$$\mathcal{L}(\boldsymbol{\omega}, \lambda) := \frac{1}{2} \left\| \mathbf{Y} \boldsymbol{\omega} \right\|_2^2 + \lambda (1 - \mathbf{1}_K^T \boldsymbol{\omega})$$

Take derivative

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\omega}} = \mathbf{Y}^T \mathbf{Y} \boldsymbol{\omega} - \lambda \mathbf{1}_K$$

and

$$\frac{\partial \mathcal{L}}{\partial \lambda} = 1 - \mathbf{1}_K^T \boldsymbol{\omega} = 0$$

Assume that \mathbf{x}_i and its K nearest vectors are linearly independent

$$\text{rank}(\mathbf{N}_i) = \text{rank}(\mathbf{N}_i^T \mathbf{N}_i) = K$$

The equation about rank is easy to prove. The column vector of \mathbf{Y} are linearly independent, if not, we can get

$$\exists \lambda_i \neq 0, i = 1, 2, \dots, K, \sum_{j=1}^K \lambda_j (\mathbf{x}_i - \mathbf{n}_j) = 0$$

and

$$\mathbf{x}_i \sum_{j=1}^K \lambda_j - \sum_{j=1}^K \lambda_j \mathbf{n}_j = 0$$

It shows that \mathbf{x}_i and its K nearest vectors will not be linearly independent which leads to a contradiction. Therefore

$$\text{rank}(\mathbf{Y}) = \text{rank}(\mathbf{Y}^T \mathbf{Y}) = K$$

If unluckily \mathbf{x}_i and its K nearest vectors not fit our assumption, we can always force our matrix to be reversible using the idea of ridge regression. For the primal problem, the optimal solution is

$$\boldsymbol{\omega}^* = (\mathbf{N}_i^T \mathbf{N}_i)^{-1} (\lambda \mathbf{1}_K + \mathbf{N}_i^T \mathbf{x}_i)$$

and for our transformed problem

$$\boldsymbol{\omega}^* = \lambda (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{1}_K$$

Let

$$\mathbf{C} := \mathbf{Y}^T \mathbf{Y} = (\mathbf{N}_i - \mathbf{x}_i \mathbf{1}_K^T)^T (\mathbf{N}_i - \mathbf{x}_i \mathbf{1}_K^T)$$

Using the constraints

$$\mathbf{1}_K^T \boldsymbol{\omega}^* = \lambda \mathbf{1}_K^T \mathbf{C}^{-1} \mathbf{1}_K = 1$$

And

$$\boldsymbol{\omega}^* = \frac{\mathbf{C}^{-1} \mathbf{1}_K}{\mathbf{1}_K^T \mathbf{C}^{-1} \mathbf{1}_K} = \text{norm}(\mathbf{C}^{-1} \mathbf{1}_K)$$

where `norm` is the function of normalization through sum all the elements of the vector

$$\text{norm}(\mathbf{x}) := \frac{1}{\sum_i x_i} \mathbf{x}$$

[2]

Notice

$$\sum_{j=1}^N w_{ij} = 1$$

and

$$\begin{aligned} \forall \mathbf{b} \in \mathbb{R}^{D \times 1}, \min_{\mathbf{z}} \sum_{i=1}^N \left\| (\mathbf{z}_i + \mathbf{b}) - \sum_{i=1}^N w_{ij} (\mathbf{z}_i + \mathbf{b}) \right\|_2^2 \\ = \min_{\mathbf{z}} \sum_{i=1}^N \left\| \mathbf{z}_i - \sum_{i=1}^N w_{ij} \mathbf{z}_i \right\|_2^2 \end{aligned}$$

[3]

select another unit orthogonal basis for the orthogonal complement

$$\mathbf{x}_n := \sum_{i=1}^L \zeta_{in} \mathbf{b}_i + \sum_{j=L+1}^D \zeta_{jn} \mathbf{b}_j$$

and the two component of \mathbf{x}_n can be defined as

$$\mathbf{x}_n^{(1)} := \sum_{i=1}^L \zeta_{in} \mathbf{b}_i$$

and

$$\mathbf{x}_n^{(2)} := \sum_{j=L+1}^D \zeta_{jn} \mathbf{b}_j$$

Notice

$$z_{in} = \langle \tilde{\mathbf{x}}_n, \mathbf{b}_i \rangle = \tilde{\mathbf{x}}_n^T \mathbf{b}_i = \mathbf{b}_i^T \tilde{\mathbf{x}}_n$$

and

$$\zeta_{in} = \langle \mathbf{x}_n, \mathbf{b}_i \rangle = \mathbf{x}_n^T \mathbf{b}_i = \mathbf{b}_i^T \mathbf{x}_n$$

Therefore

$$\frac{\partial \mathcal{L}}{\partial z_{in}} = \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{x}}_n} \frac{\partial \tilde{\mathbf{x}}_n}{\partial z_{in}} = 2(\tilde{\mathbf{x}}_n - \mathbf{x}_n)^T \mathbf{b}_i = 2(z_{in} - \zeta_{in})$$

Let the value of derivative function to zero and the optimal solution for dimension reduction satisfies

$$\begin{aligned} \tilde{\mathbf{x}}_n &= \sum_{i=1}^L z_{in} \mathbf{b}_i = \sum_{i=1}^L \zeta_{in} \mathbf{b}_i \\ &= \mathbf{x}_n^{(1)} = \pi_{\text{span}(\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_L)}(\mathbf{x}_n) \\ &= \left(\sum_{i=1}^L \mathbf{b}_i \mathbf{b}_i^T \right) \mathbf{x}_n = \mathbf{B} \mathbf{B}^T \mathbf{x}_n \end{aligned}$$

and

$$\tilde{\mathbf{X}} = \mathbf{B} \mathbf{B}^T \mathbf{X}$$

and optimal target is

$$\begin{aligned}
\min_{\mathbf{B}} \sum_{i=1}^N \left\| \mathbf{x}_i^{(2)} \right\|_2^2 &= \min_{\mathbf{B}} \sum_{i=1}^N \sum_{j=L+1}^D \zeta_{ji}^2 = \min_{\mathbf{B}} \sum_{i=1}^N \sum_{j=L+1}^D \mathbf{b}_j^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{b}_j \\
&= \min_{\mathbf{B}} \sum_{j=L+1}^D \mathbf{b}_j^T \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{b}_j, \text{ s.t. } \mathbf{B}^T \mathbf{B} = \mathbf{I}_{D-L}
\end{aligned}$$

corresponding to the perspective of maximize variance.

[4]

The low-dimensional representation in latent space whose Euclidean distance is close to the defined distance

$$d_{ij}^2 \approx \|\mathbf{z}_i - \mathbf{z}_j\|_2^2 = \|\mathbf{z}_i\|_2^2 + \|\mathbf{z}_j\|_2^2 - 2 \langle \mathbf{z}_i, \mathbf{z}_j \rangle$$

Define Gram matrix of \mathbf{Z}

$$\mathbf{G} := \{\mathbf{z}_n^T \mathbf{z}_m\}_{n,m=1}^N = \mathbf{Z}^T \mathbf{Z}$$

and we can obtain

$$g_{ij} \approx -\frac{1}{2}(d_{ij}^2 - \|\mathbf{z}_i\|_2^2 - \|\mathbf{z}_j\|_2^2) = -\frac{1}{2}(d_{ij}^2 - g_{ii} - g_{jj})$$

To simplify the problem, in classical MDS, we often assume that \mathbf{Z} has zero mean, i.e.

$$\mathbf{Z} \mathbf{1}_N = \mathbf{0}_N$$

And we also constraint the rows and columns of distance matrix to zero mean, define the margin min

$$\begin{aligned}
\bar{d}_{i\cdot}^2 &= \frac{1}{N} \sum_{j=1}^N d_{ij}^2 = \|\mathbf{z}_i\|_2^2 + \frac{1}{N} \sum_{j=1}^N \|\mathbf{z}_j\|_2^2 - \frac{2}{N} \langle \mathbf{z}_i, \mathbf{Z} \mathbf{1}_N \rangle \\
&= \|\mathbf{z}_i\|_2^2 + \frac{1}{N} \sum_{j=1}^N \|\mathbf{z}_j\|_2^2 = \|\mathbf{z}_i\|_2^2 + \frac{1}{N} \text{tr}(\mathbf{G})
\end{aligned}$$

and

$$\begin{aligned}
\bar{d}_{\cdot j}^2 &= \frac{1}{N} \sum_{i=1}^N d_{ij}^2 = \|\mathbf{z}_j\|_2^2 + \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i\|_2^2 - \frac{2}{N} \langle \mathbf{Z} \mathbf{1}_N, \mathbf{z}_j \rangle \\
&= \|\mathbf{z}_j\|_2^2 + \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i\|_2^2 = \|\mathbf{z}_j\|_2^2 + \frac{1}{N} \text{tr}(\mathbf{G})
\end{aligned}$$

and

$$\bar{d}_{\cdot\cdot}^2 = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = \frac{1}{N} \sum_{i=1}^N \bar{d}_{i\cdot}^2 = \frac{2}{N} \text{tr}(\mathbf{G})$$

so

$$\begin{aligned}
g_{nn} &= \|\mathbf{z}_n\|_2^2 \\
&= \bar{d}_{n\cdot}^2 - \frac{1}{N} \text{tr}(\mathbf{G}) = \bar{d}_{n\cdot}^2 - \frac{1}{2} \bar{d}_{\cdot\cdot}^2 \\
&= \bar{d}_{\cdot n}^2 - \frac{1}{N} \text{tr}(\mathbf{G}) = \bar{d}_{\cdot n}^2 - \frac{1}{2} \bar{d}_{\cdot\cdot}^2
\end{aligned}$$

and

$$\begin{aligned}
g_{ij} &= -\frac{1}{2} (d_{ij}^2 - g_{ii} - g_{jj}) \\
&= -\frac{1}{2} (d_{ij}^2 - \bar{d}_{i\cdot}^2 - \bar{d}_{\cdot j}^2 + \bar{d}_{\cdot\cdot}^2)
\end{aligned}$$

Using centering matrix \mathbf{C} to remove the margin mean of our defined distance, i.e. to set the sum of column or row to 0, then set the coefficient

$$\mathbf{K} = -\frac{1}{2} \mathbf{C}(\mathbf{D} \odot \mathbf{D}) \mathbf{C}$$

and \mathbf{C} is centering matrix

$$\mathbf{C} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_{N \times N}$$

and we final get

$$g_{ij} = \mathbf{z}_i^T \mathbf{z}_j \approx k_{ij}$$

Therefore, we construct the strain loss to force g_{ij} close to k_{ij}

$$\text{Strain}_d(\mathbf{Z}) = \left(\frac{\sum_{n,m=1}^N (k_{nm} - \mathbf{z}_n^T \mathbf{z}_m)^2}{\sum_{n,m=1}^N k_{nm}^2} \right)^{1/2}$$

[5]

Notice

$$\mathbf{G}^* = \mathbf{B}_L^* \mathbf{B}_L^{*T} \mathbf{K} \mathbf{B}_L^{*T} \mathbf{B}_L^*$$

It's easy to conclude that

$$\mathbf{G}^{*T} = \mathbf{G}^*$$

and

$$\mathbf{x}^T \mathbf{G} \mathbf{x} = (\mathbf{B}_L^{*T} \mathbf{B}_L^* \mathbf{x})^T \mathbf{K} (\mathbf{B}_L^{*T} \mathbf{B}_L^* \mathbf{x}) \geq 0$$

for \mathbf{K} is positive semi-definite.

[6]

It's easy to conclude that all the Φ we defined is symmetric and positive semi-definite, mostly for the reason

$$\forall \mathbf{M} \in \mathbb{R}^{n \times m}, f(\mathbf{x}) = \mathbf{x}^T \mathbf{M} \mathbf{M}^T \mathbf{x} = \|\mathbf{M}^T \mathbf{x}\|_2^2 \geq 0$$

For LLE

$$\Phi := (\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W}) \text{For constraints}$$

$$\mathbf{W} \mathbf{1}_N = \mathbf{1}_N$$

We can get

$$(\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W}) \mathbf{1}_N = (\mathbf{I} - \mathbf{W})^T (\mathbf{1}_N - \mathbf{W} \mathbf{1}_N) = \mathbf{0}_N = \mathbf{0} \mathbf{1}_N$$

Φ has a eigenvalue 0 and its corresponding eigenvector. In fact, the eigenvalue and corresponding eigenvector are always discarded for it is trivial for the problem.

For PCA

$$\Phi := \sum_{j=L+1}^D \mathbf{x}_j \mathbf{x}_j^T \in \mathbf{R}^{N \times N}$$

For Kernel PCA

$$\Phi := \sum_{j=L+1}^D \phi(\mathbf{x}_j) \phi(\mathbf{x}_j)^T \in \mathbf{R}^{N \times N}$$

For ISOMap

$$\Phi := \sum_{j=L+1}^D \mathbf{k}_j \mathbf{k}_j^T \in \mathbf{R}^{N \times N}$$

For Eigenmap

We just need to prove that for Eigenmap

$$\Phi = \text{diag}(\mathbf{A}\mathbf{1}_N) - \mathbf{A}$$

is semi-definite.

Notice

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{x}^T \Phi \mathbf{x} = \sum_{i,j=1}^n a_{ij} x_i^2 - \sum_{i,j=1}^n a_{ij} x_i x_j \\ &= \sum_{i,j=1}^n a_{ij} \left(\frac{x_i^2 + x_j^2}{2} - x_i x_j \right) \geq 0 \end{aligned}$$

It's easy to prove that Φ has a eigenvalue 0 and its corresponding eigenvector

$$\Phi \mathbf{1}_N = (\text{diag}(\mathbf{A}\mathbf{1}_N) - \mathbf{A}) \mathbf{1}_N = \mathbf{0}_N = 0 \mathbf{1}_N$$

In fact, the eigenvalue and corresponding eigenvector are always discarded for it is trivial for the problem.

[7]

Notice

$$\sum_{n=1}^D \text{tr}(\gamma_n \gamma_n^T) = \text{tr} \left(\sum_{n=1}^D \gamma_n \gamma_n^T \right) = \text{tr}(\mathbf{Y} \mathbf{Y}^T) = \text{tr}(\mathbf{Y}^T \mathbf{Y}) = L$$

and do Gram-Schmidt Orthogonalization to the column vectors of \mathbf{Y}

$$(\mathbf{Y}, \mathbf{Y}') \in \mathbb{R}^{D \times D}, (\mathbf{Y}, \mathbf{Y}')^T (\mathbf{Y}, \mathbf{Y}') = (\mathbf{Y}, \mathbf{Y}') (\mathbf{Y}, \mathbf{Y}')^T = \mathbf{I}_D$$

and

$$\mathbf{Y}'^T = (\gamma'_1, \gamma'_2, \dots, \gamma'_D), \gamma'_i \in \mathbb{R}^{(D-L) \times 1}$$

Therefore

$$\begin{aligned} (\mathbf{Y}, \mathbf{Y}') (\mathbf{Y}, \mathbf{Y}')^T &= \mathbf{I}_D \\ \Rightarrow (\gamma_i^T, \gamma_i'^T) (\gamma_i^T, \gamma_i'^T)^T &= \gamma_i^T \gamma_i + \gamma_i'^T \gamma_i' = 1 \\ \Rightarrow 0 \leq \gamma_i^T \gamma_i &= \text{tr}(\gamma_i^T \gamma_i) = \text{tr}(\gamma_i \gamma_i^T) \leq 1, i = 1, 2, \dots, D \end{aligned}$$

So

$$\sum_{n=1}^D a_n = L, 0 \leq a_n \leq 1, n = 1, 2, \dots, D$$

[8]

If not, there must exists a smaller combination

$$\exists s \geq D - L + 1, a_s < 1$$

and

$$\exists t < D - L + 1, a_t > 0$$

The solution is not the smallest solution for moving the weight and we can get

$$\sum_{n \neq s, t}^D \lambda_n a_n + \lambda_s(a_s + \varepsilon) + \lambda_t(a_t - \varepsilon) < \sum_{n=1}^D \lambda_n a_n$$

where

$$0 \leq a_s + \varepsilon \leq 1, 0 \leq a_t - \varepsilon \leq 1$$

[9]

The optimal solution satisfies

$$\Phi \mathbf{b}_i = \phi(\mathbf{X})\phi(\mathbf{X})^T \mathbf{b}_i = \lambda_i \mathbf{b}_i$$

Notice

$$\mathbf{b}_i = \frac{1}{\lambda_i} \phi(\mathbf{X})\phi(\mathbf{X})^T \mathbf{b}_i$$

Let

$$\mathbf{c}_i = \frac{1}{\lambda_i} \phi(\mathbf{X})^T \mathbf{b}_i \Rightarrow \mathbf{C} = \phi(\mathbf{X})^T \mathbf{B} \mathbf{\Lambda}^{-1}$$

and

$$\mathbf{b}_i = \phi(\mathbf{X}) \mathbf{c}_i \Rightarrow \mathbf{B} = \phi(\mathbf{X}) \mathbf{C}$$

Therefore

$$\Phi \mathbf{b}_i = \phi(\mathbf{X}) \mathbf{K} \mathbf{c}_i = \phi(\mathbf{X}) \lambda_i \mathbf{c}_i \Leftarrow \mathbf{K} \mathbf{c}_i = \lambda_i \mathbf{c}_i$$

The equation must have solution for \mathbf{K} and $\phi(\mathbf{X})\phi(\mathbf{X})^T$ have the same non-zero eigenvalue, because for $\lambda \geq 0$, we should verify that

$$\exists k \neq 0, f_{\mathbf{A}\mathbf{A}^T}(\lambda) := |\mathbf{A}\mathbf{A}^T - \lambda \mathbf{I}| = k f_{\mathbf{A}^T \mathbf{A}}(\lambda) := k |\mathbf{A}^T \mathbf{A} - \lambda \mathbf{I}|$$

Think of block matrices, assume $\mathbf{A} \in \mathbb{R}^{m \times n}$

$$\begin{pmatrix} \sigma \mathbf{I}_n & \mathbf{A}^T \\ \mathbf{A} & \sigma \mathbf{I}_m \end{pmatrix} \rightarrow \begin{pmatrix} \sigma \mathbf{I}_n & \mathbf{A}^T \\ \mathbf{O} & \sigma \mathbf{I}_m - \frac{1}{\sigma} \mathbf{A} \mathbf{A}^T \end{pmatrix}$$

和

$$\begin{pmatrix} \sigma \mathbf{I}_n & \mathbf{A}^T \\ \mathbf{A} & \sigma \mathbf{I}_m \end{pmatrix} \rightarrow \begin{pmatrix} \sigma \mathbf{I}_n - \frac{1}{\sigma} \mathbf{A}^T \mathbf{A} & \mathbf{O} \\ \mathbf{A} & \sigma \mathbf{I}_m \end{pmatrix}$$

where $\sigma \neq 0$.

For elementary transformation of partitioned matrices will not change the value of determinant, so

$$\sigma^n \left| \sigma \mathbf{I}_m - \frac{1}{\sigma} \mathbf{A} \mathbf{A}^T \right| = \sigma^m \left| \sigma \mathbf{I}_n - \frac{1}{\sigma} \mathbf{A}^T \mathbf{A} \right|$$

let $\lambda = \sigma^2$, finally we derive the result