

Manifold Learning

4.1-4.3 Introduction

20230112

Jiahao Lai

Abstract Manifold learning methods are one of the most exciting developments in machine learning in recent years. The central idea underlying these methods is that although natural data is typically represented in very high-dimensional spaces, the process generating the data is often thought to have relatively few degrees of freedom. A natural mathematical characterization of this intuition is to model the data as lying on or near a low-dimensional manifold embedded in a higher dimensional space.

In this chapter, we discuss the problem of nonlinear dimensionality reduction (NLDR): the task of recovering meaningful low-dimensional structures hidden in high-dimensional data. In many cases of interest, the observed data are found to lie on an embedded submanifold of the high-dimensional space, e.g., images generated by different views of the same three-dimensional (3D) object. Intuitively,

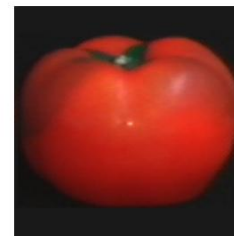
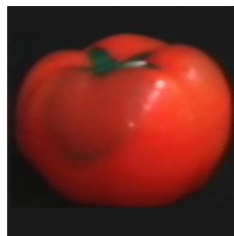
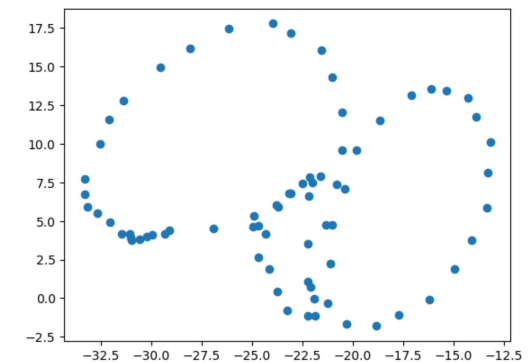


Figure 2: Tomato pictures taken from 0, 50, and 100 degree angles. In total, there are $72 (= \frac{360}{5})$ different angles.



Manifold:



Fig. 4.1 Manifold examples. Left: A curve in \mathbb{R}^3 can be considered as a 1D manifold embedded in the 3D Euclidean space. Middle: the set of points sampled from Swiss roll manifold. The set of points sampled from S-shaped manifold

Manifold: A d -dimensional manifold \mathbb{M} has the property that it is locally homeomorphic with respect to \mathbb{R}^d . That is, for each $x \in \mathbb{R}^d$, there is an open neighborhood around x , N_x , and a homeomorphism $f: N_x \rightarrow \mathbb{R}^d$. These neighborhoods are referred to as coordinate patches, and the map is referred to as a coordinate chart. The image of the coordinate chart is referred to as the parameter space.

Embedding: an embedding of a manifold \mathbb{M} into \mathbb{R}^D is a smooth homeomorphism from \mathbb{M} to a subset of \mathbb{R}^D . The algorithms discussed in this chapter find embeddings of discrete point sets, by which we simply mean a mapping of the point set into another space (typically a lower dimensional Euclidean space). An embedding of a dissimilarity matrix into \mathbb{R}^D is a configuration of points whose interpoint distances match those given in the matrix.

Graph embedding: the embedding of a graph G on a surface \mathbb{M} is a representation of G on \mathbb{M} in which the points of \mathbb{M} are associated with vertices and simple arcs (homeomorphic images of $[0,1]$) and are associated with edges in such a way that:

- the end points of the arc associated with an edge e are the points associated with the end vertices of e ;
- an arc includes no points associated with other vertices;
- two arcs never intersect at a point which is interior to either of the arcs.

Manifold learning in dimensionality reduction:

- In the classical approaches to dimensionality reduction, various methods generate linear/nonlinear maps like kernel PCA, kernel Fisher discriminant. They don't explicitly consider the structure of the manifold on which the data may possibly reside.
- Manifold learning methods address the dimensionality reduction problem by **uncovering the intrinsic low dimensional geometric structure(local or global geometry)** hidden in their high-dimensional observations and constructing a representation of the data in low-dimensional space.

How to construct a map using certain properties of the manifold?

- Global Methods: preserve metrics at all scale
MDS(preserve inner product or Euclidean distance),
Isomap(preserve geodesic distance)
- Local Methods: preserve local geometry of the data
LLE(find optimal linear reconstruction in a neighborhood),
Laplacian eigenmaps(embed vertices into a graph and use graph Laplacian to derive a smooth mapping)
- Hybrid Methods: make use of both local and global properties of the manifold
LTSA

denote $\mathbb{R}^{d \times n} \ni \mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and $\mathbb{R}^{p \times n} \ni \mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_n]$.

Classical MDS : preserves the similarity of the data.

Similarity measure of points: **inner product**

Cost function: **STRAIN**

$$\underset{\{\mathbf{y}_i\}_{i=1}^n}{\text{minimize}} \quad c_1 := \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i^\top \mathbf{x}_j - \mathbf{y}_i^\top \mathbf{y}_j)^2, \quad (1)$$

whose matrix form is:

$$\underset{\mathbf{Y}}{\text{minimize}} \quad c_1 = \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2, \quad (2)$$

where $\|\cdot\|_F$ denotes the Frobenius norm, and $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{Y}^\top \mathbf{Y}$ are the Gram matrices of the original data \mathbf{X} and the embedded data \mathbf{Y} , respectively.

The objective function, in Eq. (2), is simplified as:

$$\begin{aligned} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2 &= \text{tr}[(\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y})^\top (\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y})] \\ &= \text{tr}[(\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y})(\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y})] \\ &= \text{tr}[(\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y})^2], \end{aligned}$$

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \Delta \mathbf{V}^\top, \quad (3)$$

$$\mathbf{Y}^\top \mathbf{Y} = \mathbf{Q} \Psi \mathbf{Q}^\top, \quad (4)$$

$$\begin{aligned} \therefore \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2 &= \text{tr}[(\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y})^2] \\ &= \text{tr}[(\mathbf{V} \Delta \mathbf{V}^\top - \mathbf{Q} \Psi \mathbf{Q}^\top)^2] \\ &\stackrel{(a)}{=} \text{tr}[(\mathbf{V} \Delta \mathbf{V}^\top - \mathbf{V} \mathbf{V}^\top \mathbf{Q} \Psi \mathbf{Q}^\top \mathbf{V} \mathbf{V}^\top)^2] \\ &= \text{tr}[(\mathbf{V}(\Delta - \mathbf{V}^\top \mathbf{Q} \Psi \mathbf{Q}^\top \mathbf{V})\mathbf{V}^\top)^2] \\ &= \text{tr}[\mathbf{V}^2(\Delta - \mathbf{V}^\top \mathbf{Q} \Psi \mathbf{Q}^\top \mathbf{V})^2(\mathbf{V}^\top)^2] \\ &\stackrel{(b)}{=} \text{tr}[(\mathbf{V}^\top)^2 \mathbf{V}^2 (\Delta - \mathbf{V}^\top \mathbf{Q} \Psi \mathbf{Q}^\top \mathbf{V})^2] \\ &= \text{tr}[(\underbrace{\mathbf{V}^\top \mathbf{V}}_I)^2 (\Delta - \mathbf{V}^\top \mathbf{Q} \Psi \mathbf{Q}^\top \mathbf{V})^2] \\ &\stackrel{(c)}{=} \text{tr}[(\Delta - \mathbf{V}^\top \mathbf{Q} \Psi \mathbf{Q}^\top \mathbf{V})^2], \end{aligned}$$

Let $\mathbb{R}^{n \times n} \ni \mathbf{M} := \mathbf{V}^\top \mathbf{Q}$, so:

$$\|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2 = \text{tr}[(\Delta - \mathbf{M} \Psi \mathbf{M}^\top)^2].$$

$$\begin{aligned} c_1 &= \text{tr}[(\Delta - \mathbf{M} \Psi \mathbf{M}^\top)^2] \\ &= \text{tr}(\Delta^2 + (\mathbf{M} \Psi \mathbf{M}^\top)^2 - 2\Delta \mathbf{M} \Psi \mathbf{M}^\top) \\ &= \text{tr}(\Delta^2) + \text{tr}((\mathbf{M} \Psi \mathbf{M}^\top)^2) - 2\text{tr}(\Delta \mathbf{M} \Psi \mathbf{M}^\top). \end{aligned}$$

$$\begin{aligned}
\mathbb{R}^{n \times n} \ni \frac{\partial c_1}{\partial \mathbf{M}} &= 2(\mathbf{M}\Psi\mathbf{M}^\top)\mathbf{M}\Psi - 2\Delta\mathbf{M}\Psi \stackrel{\text{set}}{=} \mathbf{0} \\
\implies (\mathbf{M}\Psi\mathbf{M}^\top)(\mathbf{M}\Psi) &= (\Delta)(\mathbf{M}\Psi) \\
\stackrel{(a)}{\implies} \mathbf{M}\Psi\mathbf{M}^\top &= \Delta,
\end{aligned} \tag{5}$$

where (a) is because $\mathbf{M}\Psi \neq \mathbf{0}$.

$$\begin{aligned}
\mathbb{R}^{n \times n} \ni \frac{\partial c_1}{\partial \Psi} &= 2\mathbf{M}^\top(\mathbf{M}\Psi\mathbf{M}^\top)\mathbf{M} - 2\mathbf{M}^\top\Delta\mathbf{M} \stackrel{\text{set}}{=} \mathbf{0} \\
\implies \mathbf{M}^\top(\mathbf{M}\Psi\mathbf{M}^\top)\mathbf{M} &= \mathbf{M}^\top(\Delta)\mathbf{M} \\
\stackrel{(a)}{\implies} \mathbf{M}\Psi\mathbf{M}^\top &= \Delta,
\end{aligned} \tag{6}$$

whose one possible solution is:

$$\mathbf{M} = \mathbf{I}, \tag{7}$$

$$\Psi = \Delta. \tag{8}$$

which means that the minimum value of the non-negative objective function $\text{tr}((\Delta - \mathbf{M}\Psi\mathbf{M}^\top)^2)$ is zero.

We had $\mathbf{M} = \mathbf{V}^\top \mathbf{Q}$. Therefore, according to Eq. (7), we have:

$$\therefore \mathbf{V}^\top \mathbf{Q} = \mathbf{I} \implies \mathbf{Q} = \mathbf{V}. \tag{9}$$

According to Eq. (4), we have:

$$\begin{aligned}
\mathbf{Y}^\top \mathbf{Y} &= \mathbf{Q}\Psi\mathbf{Q}^\top \stackrel{(a)}{=} \mathbf{Q}\Psi^{\frac{1}{2}}\Psi^{\frac{1}{2}}\mathbf{Q}^\top \implies \mathbf{Y} = \Psi^{\frac{1}{2}}\mathbf{Q}^\top \\
&\stackrel{(8),(9)}{\implies} \mathbf{Y} = \Delta^{\frac{1}{2}}\mathbf{V}^\top,
\end{aligned} \tag{10}$$

In summary, for embedding \mathbf{X} using classical MDS, the eigenvalue decomposition of $\mathbf{X}^\top \mathbf{X}$ is obtained as in Eq. (3). Then, using Eq. (10), $\mathbf{Y} \in \mathbb{R}^{n \times n}$ is obtained. Truncating this \mathbf{Y} to have $\mathbf{Y} \in \mathbb{R}^{p \times n}$, with the first (top) p rows, gives us the p -dimensional embedding of the n points. Note that the leading p columns are used because singular values are sorted from largest to smallest in SVD which can be used for Eq. (3).

Kernel MDS / generalized classical MDS:

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

where $\mathbb{R}^{n \times n} \ni \mathbf{G} := \mathbf{X}^\top \mathbf{X}$ is the Gram matrix. If $\mathbb{R}^n \ni \mathbf{g} := [\mathbf{g}_1, \dots, \mathbf{g}_n] = [\mathbf{G}_{11}, \dots, \mathbf{G}_{nn}] = \text{diag}(\mathbf{G})$, we have:

$$\begin{aligned}
d_{ij}^2 &= \mathbf{g}_i - 2\mathbf{G}_{ij} + \mathbf{g}_j, \\
\mathbf{D} &= \mathbf{g}\mathbf{1}^\top - 2\mathbf{G} + \mathbf{1}\mathbf{g}^\top = \mathbf{1}\mathbf{g}^\top - 2\mathbf{G} + \mathbf{g}\mathbf{1}^\top,
\end{aligned}$$

where $\mathbf{1}$ is the vector of ones and \mathbf{D} is the distance matrix with squared Euclidean distance (d_{ij}^2 as its elements). Let $\mathbb{R}^{n \times n} \ni \mathbf{H} := \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$ denote the centering matrix. We

double-center the matrix D as follows (Oldford, 2018):

$$\begin{aligned}
HDH &= (I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)D(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top) \\
&= (I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)(\mathbf{1}\mathbf{g}^\top - 2G + \mathbf{g}\mathbf{1}^\top)(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top) \\
&= \underbrace{[(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)\mathbf{1}\mathbf{g}^\top - 2(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)G]}_{=0} \\
&\quad + (I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)\mathbf{g}\mathbf{1}^\top](I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top) \\
&= -2(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)G(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top) \\
&\quad + (I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)\mathbf{g}\mathbf{1}^\top \underbrace{(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)}_{=0} \\
&= -2(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)G(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top) = -2HGH
\end{aligned}$$

$$\therefore HGH = HX^\top XH = -\frac{1}{2}HDH. \quad (11)$$

If data X are already centered, i.e., the mean has been removed ($X \leftarrow XH$), Eq. (11) becomes:

$$X^\top X = -\frac{1}{2}HDH. \quad (12)$$

use kernel matrix rather than gram matrix:

$$\mathbb{R}^{n \times n} \ni K = -\frac{1}{2}HDH. \quad (13)$$

$$K = V\Delta V^\top. \quad (14)$$

Then, using Eq. (10), $Y \in \mathbb{R}^{n \times n}$ is obtained. It is noteworthy that in this case, we are replacing $X^\top X$ with the kernel $K = \Phi(X)^\top \Phi(X)$ and then, according to Eqs. (10) and (14), we have:

$$K = Y^\top Y. \quad (15)$$

Truncating the Y , obtained from Eq. (10), to have $Y \in \mathbb{R}^{p \times n}$, with the first (top) p rows, gives us the p -dimensional embedding of the n points. It is noteworthy that, because of using kernel in the generalized classical MDS, one can name it the *kernel classical MDS*.

Metric MDS : preserves the distance of points rather than the similarity.

Cost function: **stress/sstress function**

$$\begin{aligned} & \text{minimize}_{\{\mathbf{y}_i\}_{i=1}^n} \\ c_2 &:= \left(\frac{\sum_{i=1}^n \sum_{j=1, j < i}^n (d_x(\mathbf{x}_i, \mathbf{x}_j) - d_y(\mathbf{y}_i, \mathbf{y}_j))^2}{\sum_{i=1}^n \sum_{j=1, j < i}^n d_x^2(\mathbf{x}_i, \mathbf{x}_j)} \right)^{\frac{1}{2}}, \end{aligned} \quad (16)$$

or, without the normalization factor:

$$\begin{aligned} & \text{minimize}_{\{\mathbf{y}_i\}_{i=1}^n} \\ c_2 &:= \left(\sum_{i=1}^n \sum_{j=1, j < i}^n (d_x(\mathbf{x}_i, \mathbf{x}_j) - d_y(\mathbf{y}_i, \mathbf{y}_j))^2 \right)^{\frac{1}{2}}, \end{aligned} \quad (17)$$

where $d_x(., .)$ and $d_y(., .)$ denote the distance metrics in the input and the embedded spaces, respectively.

the matrix form: $\|\mathbf{D} - \mathbf{D}'\|_F^2$

Other criteria for MDS have also been studied, for example, the STRESS and SSTRESS given by $\|\mathbf{D} - \mathbf{D}'\|_F^2$ and $\|\mathbf{D}^2 - \mathbf{D}'^2\|_F^2$, respectively. However, the

Use gradient descent or Newton's method to solve (16) since it doesn't have closed-form solution.

Non-Metric MDS :

In *non-metric MDS*, rather than using a distance metric, $d_y(\mathbf{x}_i, \mathbf{x}_j)$, for the distances between points in the embedding space, we use $f(d_y(\mathbf{x}_i, \mathbf{x}_j))$ where $f(.)$ is a non-parametric monotonic function. In other words, only the order of dissimilarities is important rather than the amount of dissimilarities (Agarwal et al., 2007; Jung, 2013):

$$\begin{aligned} d_y(\mathbf{y}_i, \mathbf{y}_j) \leq d_y(\mathbf{y}_k, \mathbf{y}_\ell) &\iff \\ f(d_y(\mathbf{y}_i, \mathbf{y}_j)) &\leq f(d_y(\mathbf{y}_k, \mathbf{y}_\ell)). \end{aligned} \quad (20)$$

The optimization in non-metric MDS is (Agarwal et al., 2007):

$$\begin{aligned} & \text{minimize}_{\{\mathbf{y}_i\}_{i=1}^n} \quad c_3 := \\ & \left(\frac{\sum_{i=1}^n \sum_{j=1, j < i}^n (d_x(\mathbf{x}_i, \mathbf{x}_j) - f(d_y(\mathbf{y}_i, \mathbf{y}_j)))^2}{\sum_{i=1}^n \sum_{j=1, j < i}^n d_x^2(\mathbf{x}_i, \mathbf{x}_j)} \right)^{\frac{1}{2}}. \end{aligned} \quad (21)$$

Isometric Feature Mapping(Isomap):

In classical MDS, we use Euclidean distance

$$\underset{\{\mathbf{y}_i\}_{i=1}^n}{\text{minimize}} \quad c_1 := \sum_{i=1}^n \sum_{j=1}^n (\mathbf{x}_i^\top \mathbf{x}_j - \mathbf{y}_i^\top \mathbf{y}_j)^2,$$

$$\mathbf{X}^\top \mathbf{X} = -\frac{1}{2} \mathbf{H} \mathbf{D} \mathbf{H}.$$

which is not suitable for a nonlinearly embedded manifold. Thus in kernel MDS, we use kernel matrix to substitute the gram matrix.

$$\mathbb{R}^{n \times n} \ni \mathbf{K} = -\frac{1}{2} \mathbf{H} \mathbf{D} \mathbf{H}.$$

Isomap is a special case of kernel MDS. Isomap uses an approximation of the geodesic distance in above matrix \mathbf{D} .

Geodesic distance: the length of shortest path between two points on the possibly curvy manifold.



Calculation of the geodesic distance is very difficult because it requires the knowledge of differential geometry. Therefore, Isomap approximates the geodesic distance by piecewise Euclidean distances.

The shortest path is found using algorithms such as the Dijkstra algorithm or the Floyd-Warshal algorithm.

$$D_{ij}^{(g)} := \min_{\mathbf{r}} \sum_{i=2}^l \|\mathbf{r}_i - \mathbf{r}_{i+1}\|_2,$$

$$\mathbb{R}^{n \times n} \ni \mathbf{K} = -\frac{1}{2} \mathbf{H} \mathbf{D}^{(g)} \mathbf{H}.$$

1. Find the neighborhood for each point \mathbf{x}_i .
2. Compute the shortest path distance between all pairs of points using Dijkstra's or Floyd's algorithm, and store the squares of these distances in a matrix.
3. Apply classical MDS to the distance matrix.

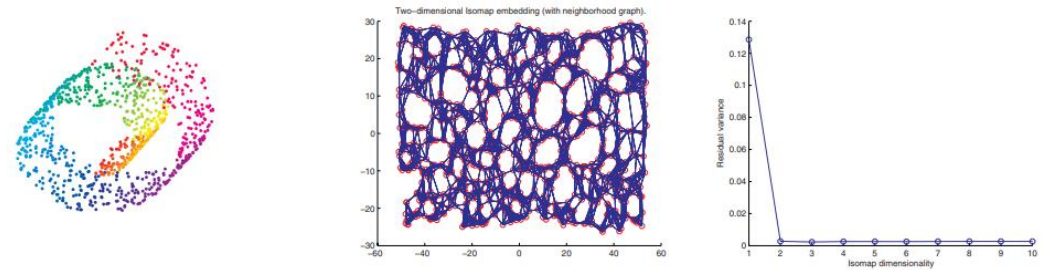


Fig. 4.4 The Swiss roll data set, illustrating the 2D embedding recovered by the Isomap in step three, and how the residual variance decreases as the dimensionality d is increased. Left: the true manifold; middle: the 2D embedding recovered by the Isomap with neighborhood graph overlaid; right: the residual variance decreases as dimensionality increases

Landmark Isomap:

One drawback of Isomap is its quadratic memory requirement: the geodesic distance matrix is dense, making the Isomap infeasible for large data set.

Nystrom Approximation:

Nystrom approximation is a technique used to approximate a positive semi-definite matrix using merely a subset of its columns (or rows) (Williams & Seeger, 2001). Consider a positive semi-definite matrix $\mathbb{R}^{n \times n} \ni \mathbf{K} \succeq 0$ whose parts are:

$$\mathbb{R}^{n \times n} \ni \mathbf{K} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{B}^\top & \mathbf{C} \end{array} \right], \quad (52)$$

where $\mathbf{A} \in \mathbb{R}^{m \times m}$, $\mathbf{B} \in \mathbb{R}^{m \times (n-m)}$, and $\mathbf{C} \in \mathbb{R}^{(n-m) \times (n-m)}$ in which $m \ll n$.

\mathbf{K} is the similarity matrix or distance matrix. Instead of computing \mathbf{K} , we designate m data points to be landmark points and compute \mathbf{A} and \mathbf{B} . Our goal is to approximate \mathbf{C} .

As the matrix \mathbf{K} is positive semi-definite, by definition, it can be written as $\mathbf{K} = \mathbf{O}^\top \mathbf{O}$. If we take $\mathbf{O} = [\mathbf{R}, \mathbf{S}]$ where \mathbf{R} are the selected columns (landmarks) of \mathbf{O} and \mathbf{S} are the other columns of \mathbf{O} . We have:

$$\mathbf{K} = \mathbf{O}^\top \mathbf{O} = \begin{bmatrix} \mathbf{R}^\top \\ \mathbf{S}^\top \end{bmatrix} [\mathbf{R}, \mathbf{S}] \quad (53)$$

$$= \begin{bmatrix} \mathbf{R}^\top \mathbf{R} & \mathbf{R}^\top \mathbf{S} \\ \mathbf{S}^\top \mathbf{R} & \mathbf{S}^\top \mathbf{S} \end{bmatrix} \stackrel{(52)}{=} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix}. \quad (54)$$

Hence, we have $\mathbf{A} = \mathbf{R}^\top \mathbf{R}$. The eigenvalue decomposition (Ghojogh et al., 2019a) of \mathbf{A} gives:

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top \quad (55)$$

$$\Rightarrow \mathbf{R}^\top \mathbf{R} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top \Rightarrow \mathbf{R} = \mathbf{\Sigma}^{(1/2)} \mathbf{U}^\top. \quad (56)$$

Moreover, we have $\mathbf{B} = \mathbf{R}^\top \mathbf{S}$ so we have:

$$\begin{aligned} \mathbf{B} &= (\mathbf{\Sigma}^{(1/2)} \mathbf{U}^\top)^\top \mathbf{S} = \mathbf{U} \mathbf{\Sigma}^{(1/2)} \mathbf{S} \\ \stackrel{(a)}{\Rightarrow} \mathbf{U}^\top \mathbf{B} &= \mathbf{\Sigma}^{(1/2)} \mathbf{S} \Rightarrow \mathbf{S} = \mathbf{\Sigma}^{(-1/2)} \mathbf{U}^\top \mathbf{B}, \end{aligned} \quad (57)$$

$$\begin{aligned} \mathbf{C} &= \mathbf{S}^\top \mathbf{S} = \mathbf{B}^\top \mathbf{U} \mathbf{\Sigma}^{(-1/2)} \mathbf{\Sigma}^{(-1/2)} \mathbf{U}^\top \mathbf{B} \\ &= \mathbf{B}^\top \mathbf{U} \mathbf{\Sigma}^{-1} \mathbf{U}^\top \mathbf{B} \stackrel{(55)}{=} \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{B}. \end{aligned} \quad (58)$$

Therefore, Eq. (52) becomes:

$$\mathbf{K} \approx \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{B}^\top & \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{B} \end{array} \right]. \quad (59)$$

After computing \mathbf{K} , we need to do eigenvalue decomposition of the kernel matrix and compute \mathbf{Y} .

$$\mathbf{K} = \mathbf{V} \mathbf{\Delta} \mathbf{V}^\top, \quad \mathbf{K} = \mathbf{Y}^\top \mathbf{Y}, \quad \mathbf{Y} = \mathbf{\Delta}^{\frac{1}{2}} \mathbf{V}^\top$$

Recall that Eq. (14) decomposes the kernel matrix into eigenvectors and then Eq. (10) embeds data. However, for big data, the eigenvalue decomposition of kernel matrix is intractable. Therefore, using Eq. (55), we decompose an $m \times m$ submatrix of kernel. Comparing Eqs. (15) and (53) shows that:

$$\mathbb{R}^{n \times n} \ni \mathbf{Y} = [\mathbf{R}, \mathbf{S}] \stackrel{(a)}{=} [\mathbf{\Sigma}^{(1/2)} \mathbf{U}^\top, \mathbf{\Sigma}^{(-1/2)} \mathbf{U}^\top \mathbf{B}], \quad (60)$$