

# Nonlinear Dimensionality Reduction



Xian-Hong Wang (王先弘)

Department of Mathematics, National Central University  
Jhongli District, Taoyuan City 32001, Taiwan

Revised on March 2, 2023

# Outline

---

- Linear dimension reduction
  - Principal component analysis (PCA)
  - Linear discriminant analysis (LDA)
- Non-linear dimension reduction (Manifold Learning)
  - Mutidimensional Scaling (MDA)
  - Isomap
  - Diffusion Maps
  - Laplacian Eigenmaps (LE)
  - Locally linear embedding (LLE)
- Kernel method (KPCA)

## Principal component analysis (PCA)

## The basic idea of PCA

---

Given a centered dataset

$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_n] \in \mathbb{R}^{m \times n}.$$

PCA want to find a direction  $\mathbf{q}$  which maximizes the spread(variance).

Representing  $\mathbf{X}$  in terms of its thin-SVD

$$\mathbf{X} = \mathbf{U}_r \mathbf{D}_r \mathbf{V}_r^\top,$$

where  $\mathbf{U}_r \in \mathbb{R}^{n \times r}$ ,  $\mathbf{V}_r \in \mathbb{R}^{p \times r}$  are matrices with orthonormal columns,

$$\mathbf{U}_r = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_r], \mathbf{V}_r = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_r],$$

and  $\mathbf{D}_r \in \mathbb{R}^{r \times r}$  is a diagonal matrix with the nonzero singular values as diagonal entries.

## Spread

---

Let  $\mathbf{q} \in \mathbb{R}^m$  be a **unit vector**. The projections of the data vectors along the direction  $\mathbf{q}$  are given by components of the row vector

$$\mathbf{Y} = \mathbf{q}^\top \mathbf{X} = [\mathbf{q}^\top \mathbf{x}_1 \quad \mathbf{q}^\top \mathbf{x}_2 \quad \cdots \quad \mathbf{q}^\top \mathbf{x}_n] = [y_1 \quad y_2 \quad \cdots \quad y_n] \in \mathbb{R}^{1 \times n}.$$

Define the *spread* of the components as

$$\text{spread}(\mathbf{Y}) = \left( \sum_{j=1}^p y_j^2 \right)^{1/2},$$

where

$$\sum_{j=1}^N y_j^2 = \mathbf{Y} \mathbf{Y}^\top = \mathbf{q}^\top \mathbf{X} \mathbf{X}^\top \mathbf{q} = \|\mathbf{X}^\top \mathbf{q}\|_2^2.$$

## Optimization problem

---

So the optimization problem is

$$\max\{spread(\mathbf{Y})\} = \max\{\|\mathbf{X}^\top \mathbf{q}\| \mid \|\mathbf{q}\| = 1\}.$$

Therefore, it follows by the definition of **induced matrix norm** that

$$\max\{\|\mathbf{X}^\top \mathbf{q}\| \mid \|\mathbf{q}\| = 1\} = \|\mathbf{X}^\top\|_2.$$

Furthermore,

$$\max\{spread(\mathbf{Y})\} = \sigma_1 = \|\mathbf{X}^\top \mathbf{u}_1\|_2,$$

where  $\sigma_1$  is first singular value of  $\mathbf{X}$ , and  $\mathbf{u}_1$  is first left singular vector corresponding to the  $\mathbf{u}_1$ .

So

$$\mathbf{q} = \mathbf{u}_1,$$

and the maximum spread equals the largest singular value of the data matrix.

## Visualization

Subtract from the data matrix the rank 1 matrix corresponding to the first singular triplet,

$$\tilde{\mathbf{X}} = \mathbf{X} - \sigma_1 \mathbf{u}_1 (\mathbf{v}_1)^\top = \sum_{j=2}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^\top,$$

Define the spread of the components of

$$\tilde{\mathbf{Y}} = \mathbf{q}^\top \tilde{\mathbf{X}}.$$

Maximizing the spread with respect to the direction vector  $\mathbf{q}$  leads to the formula

$$\max\{\text{spread}(\tilde{\mathbf{Y}})\} = \sigma_2 = \|\tilde{\mathbf{X}}^\top \mathbf{u}_2\|,$$

confirming that the second best projection direction is  $\mathbf{u}_2$ .

## References

---

- D. Calvetti and E. Somersalo, *Mathematics of Data Science: A Computational Approach to Clustering and Classification*, SIAM, 2021, pp. 21-32.
- C. Heeyoul and C. Seungjin. Robust kernel Isomap, *ScienceDirect*, 2006, pp. 853-862.



## Linear discriminant analysis

The data considered in this section are assumed to be labeled, or annotated, in the sense that each datapoint belongs to one and only one class, and that the class-belonging information is available. More formally, let  $\mathcal{D}$  denoted the data set partitioned into  $k$  nonintersecting subsets,

$$\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \cdots \cup \mathcal{D}_k, \quad \mathcal{D}_i \cap \mathcal{D}_j = \emptyset \text{ if } i \neq j,$$

and let

$$I_j = \ell \text{ equivalent to } x^{(j)} \in \mathcal{D}_\ell, \quad 1 \leq j \leq n,$$

define the *annotation vector*  $I \in \{1, 2, \dots, k\}^j$  that contains the labels assigning each data vector to one and only one subset.

## LDA

Given a data matrix  $X \in \mathbb{R}^{m \times n}$ , center or noncentered, the components of the data vectors  $x^{(j)}$  on a given direction vector  $q \in \mathbb{R}^m$ , with  $\|q\| = 1$ , can be collected into a row vector  $y$  of length  $n$ ,

$$y = q^T X.$$

Recalling that the square of the spread of the data in the direction  $q$  is the quantity

$$\text{spread}(y)^2 = \sum_{j=1}^n y_j^2 = yy^T,$$

and replacing  $y$  by its expression in terms of the data matrix, we obtain

$$\text{spread}(y)^2 = q^T X X^T q = q^T S q.$$

The matrix  $S$ ,

$$S = X X^T \in \mathbb{R}^{m \times m},$$

is called the *scatter matrix* of the data  $X$ .

Therefore, the direction of maximum (minimum) spread of the data is the vector that maximizes (minimizes) the quadratic expression. The observation motivates us to define a family of scatter matrices to help us optimize the search for projection directions.

Without loss of generality we can reorder the data so that the columns corresponding to points in the same cluster are adjacent, i.e.,

$$\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2 \quad \cdots \quad \mathbf{X}_k] \in \mathbb{R}^{m \times n},$$

where the columns of  $\mathbf{X}_\ell$  are the data vectors in cluster  $\mathcal{D}_\ell$ ,

$$\mathbf{X}_\ell = [x^{(j_1)} \quad x^{(j_2)} \quad \cdots \quad x^{(j_{n_\ell})}] \in \mathbb{R}^{m \times n_\ell},$$

and

$$I_\ell = [j_1 \quad j_2 \quad \cdots \quad j_{n_\ell}].$$

We have

$$n_1 + n_2 + \cdots + n_k = n.$$

We define the *cluster centroids* as

$$c^{(\ell)} = \frac{1}{n_\ell} \sum_{j \in I_\ell} x^{(j)}, \quad 1 \leq \ell \leq k.$$

The *global centroid* of the data is defined by the formula

$$c = \frac{1}{n} \sum_{j=1}^n x^{(j)}.$$

We start by defining the centered data matrices for each cluster,

$$\mathbf{X}_{\ell,c} = \begin{bmatrix} x^{(j_1)} - c^{(\ell)} & x^{(j_2)} - c^{(\ell)} & \cdots & x^{(j_{n_\ell})} - c^{(\ell)} \end{bmatrix} \in \mathbb{R}^{m \times n_\ell}, \quad 1 \leq \ell \leq k.$$

The centered matrices are then collected in the *within-cluster centered data matrix*,

$$\mathbf{X}_w = [\mathbf{X}_{1,c} \quad \mathbf{X}_{2,c} \quad \cdots \quad \mathbf{X}_{k,c}] \in \mathbb{R}^{m \times n}.$$

Geometrically,  $\mathbf{X}_w$  provides a representation of the data where each cluster has been separately centered around the origin.

We define the *within-cluster scatter matrix*  $\mathbf{S}_w \in \mathbb{R}^{m \times m}$  for the centered data  $\mathbf{X}_w$  by the formula

$$\mathbf{S}_w = \mathbf{X}_w \mathbf{X}_w^\top = \sum_{\ell=1}^k \mathbf{S}_\ell$$

where  $\mathbf{S}_\ell$  is the scatter matrix of the  $\ell$ th cluster,

$$\mathbf{S}_\ell = \mathbf{X}_{\ell,c} \mathbf{X}_{\ell,c}^\top = \sum_{j \in I_\ell} (x^{(j)} - c^{(\ell)})(x^{(j)} - c^{(\ell)})^\top \in \mathbb{R}^{m \times m}, \quad 1 \leq \ell \leq k.$$

We need measure for **how the clusters themselves are scattered**. Therefore, we define a centroid approximation of the data matrix  $X$  by simply replacing each data vector  $x^{(j)}$  by the centroid of its cluster. Denote by  $\bar{X}_\ell$  the centroid approximation of  $X_\ell$ ,

$$\bar{X}_\ell = [c^{(\ell)} \quad c^{(\ell)} \quad \cdots \quad c^{(\ell)}] \in \mathbb{R}^{m \times n_\ell},$$

which is the rank 1 matrix obtained by replacing each point in the  $\ell$ th cluster by the cluster centroid. We collect the centroid approximation matrices for all the clusters into the matrix

$$\bar{X} = [\bar{X}_1 \quad \bar{X}_2 \quad \cdots \quad \bar{X}_k] \in \mathbb{R}^{m \times n},$$

whose rank is at most  $k$ .

## LDA

Next we center the matrix  $\bar{\mathbf{X}}$  by subtracting the global centroid from each column,

$$\bar{\mathbf{X}}_c = \bar{\mathbf{X}} - [\mathbf{c} \quad \mathbf{c} \quad \cdots \quad \mathbf{c}],$$

and define the *between-cluster scatter matrix*  $\mathbf{S}_b \in \mathbb{R}^{m \times m}$  by

$$\mathbf{S}_b = \bar{\mathbf{X}}_c \bar{\mathbf{X}}_c^\top,$$

which can be expressed also in the form

$$\begin{aligned} \mathbf{S}_b &= \sum_{\ell=1}^k \sum_{j=1}^{n_\ell} (\mathbf{c}^{(\ell)} - \mathbf{c})(\mathbf{c}^{(\ell)} - \mathbf{c})^\top \\ &= \sum_{\ell=1}^k n_\ell (\mathbf{c}^{(\ell)} - \mathbf{c})(\mathbf{c}^{(\ell)} - \mathbf{c})^\top. \end{aligned}$$

We have now introduced all the tools needed to formally state and pursue the objective of the LDA of finding optimal projection directions along which clusters are maximally separated.



The goal of this section is to find few projection directions in the data space along which the within-cluster spread is as small as possible and between-cluster spread as large as possible.

The objective function guiding the search for vectors determining the directions with the desired properties is the ratio

$$H(q) = \frac{q^T S_b q}{q^T S_w q}.$$

We argue that any nonzero vector  $q$  that maximizes this ratio is a desired direction.

For the time being, we assume that the matrix  $S_w$  is symmetric positive definite. The case of symmetric positive semidefinite  $S_w$  will be addressed later.

We start by observing that since for ever  $\alpha > 0$ ,

$$H(\alpha q) = H(q),$$

In particular, we can scale  $q$  so as to guarantee that

$$q^T S_w q = 1,$$

so that the objective function becomes

$$H(q) = \frac{q^T S_b q}{q^T S_w q} = q^T S_b q.$$

This observation allows us to restate the problem of finding the optimal direction in the following way.

## LDA

Since, by assumption, the matrix  $S_w$  is symmetric positive definite, it admits the Cholesky factorization,

$$S_w = K^T K,$$

where  $K$  is upper triangular and invertible. Hence we can write the constraint as

$$q^T K^T K q = \|Kq\|^2 = 1,$$

or, by defining the vector  $w$  of unit length by

$$Kq = w, \text{ or } q = K^{-1}w,$$

and substituting the expression of  $q$  in terms of  $w$ , the problem can be restated as

$$\arg \max_w w^T K^{-T} S_b K^{-1} w \quad \text{subject to } \|w\|^2 = 1,$$

where  $K^{-T} = (K^{-1})^T$ .

Following the method of Lagrange multipliers, we define a Lagrange function

$$L_{\lambda}(w) = w^{\top} K^{-\top} S_b K^{-1} w - \lambda (\|w\|^2 - 1),$$

where  $\lambda \in \mathbb{R}$  is a Lagrange multiplier, and look for its critical points, i.e., vectors  $w$  such that

$$\nabla_w L_{\lambda}(w) = 2(K^{-\top} S_b K^{-1} w - \lambda w) = 0.$$

Equivalently, the critical points must satisfy

$$K^{-\top} S_b K^{-1} w = \lambda w,$$

that is,  $w$  is an eigenvector of unit length of the matrix  $K^{-\top} S_b K^{-1}$ . Furthermore, since for such eigenvectors,

$$w^{\top} K^{-\top} S_b K^{-1} w = \lambda \|w\|^2 = \lambda,$$

the maximizer of the objective function is an eigenvector corresponding to the largest eigenvalue of the matrix, denoted by  $\lambda_1$ .

## LDA

Having found the vector  $w$ , we may solve for  $q$ . Observe that from solution, it follows that, by multiplying both sides by the matrix  $K^{-1}$ , we have

$$K^{-1}K^{-T}S_bK^{-1}w = \lambda K^{-1}w,$$

or

$$S_w^{-1}S_bq = \lambda q,$$

that is, the vectors  $q$  are eigenvectors of the nonsymmetric matrix  $S_w^{-1}S_b$ , which has thus been shown to have all real eigenvalues.

If there are more than two clusters in the data, it may be useful to find more than just one projection direction to represent the data. A natural way is to find the eigenvectors  $v^{(\ell)}$ ,  $\ell > 1$ , set  $w^{(\ell)} = v^{(\ell)}$ , and define the projection directions as

$$q^{(\ell)} = K^{-1}w^{(\ell)}.$$

Observe that while the vectors  $v^{(\ell)}$  are mutually orthogonal, the vectors  $v^{(\ell)}$  are usually not.

So far, we have assumed that the matrix  $S_w$  is symmetric positive definite, and therefore invertible. If the assumption is not valid, the argument above does not work, as  $S_w$  will not admit the Cholesky factorization. On the other hand, being symmetric positive semidefinite, the matrix  $S_w$ , could be factored as

$$S_w = S_w^{1/2} S_w^{1/2},$$

but since the factors would not be invertible, the argument does not work. To overcome this problem, we replace  $S_w$  by a nearly positive definite approximation by setting

$$S_{w,\epsilon} = S_w + \epsilon I_n,$$

where  $\epsilon > 0$  is a small scalar. For any  $v \in \mathbb{R}^m, v \neq 0$ , we have

$$v^\top S_{w,\epsilon} v = v^\top S_w v + \epsilon \|v\|^2 \geq \epsilon \|v\|^2 > 0,$$

so the existence of the Cholesky decomposition is warranted.

To understand under which conditions  $S_w$  may not be invertible, assume that  $m < n$  and consider the singular value decomposition of the within-cluster centered data,

$$X_w = UDV^\top.$$

If the rank  $r$  of the matrix  $X_w$  is less than  $n$ , the within-cluster centered data are restricted to an  $r$ -dimensional subspace of  $\mathbb{R}^n$ . On the other hand,

$$S_w = X_w X_w^\top = U D D^\top U^\top = U W U^\top,$$

with

$$W = D D^\top = \begin{bmatrix} d_1^2 & & \\ & \ddots & \\ & & d_m^2 \end{bmatrix}.$$

where  $d_j$  denotes the  $j$ th singular value of  $X_w$

Since the eigenvalues of  $S_w$  are  $d_j^2$ , the matrix  $S_w$  has only  $r$  nonzero eigenvalues, and therefore cannot be invertible, and the argument above falls apart. Effectively, the approximation is tantamount to perturbing the within-cluster centered data off the  $r$ -dimensional subspace.

Observe that the perturbed scatter matrix can be decomposed as

$$S_w + \epsilon I_m = U W U^\top + \epsilon U U^\top = U W_\epsilon U^\top,$$

where

$$W_\epsilon = \begin{bmatrix} d_1^2 + \epsilon & & \\ & \ddots & \\ & & d_m^2 + \epsilon \end{bmatrix},$$

and since all its eigenvalue are strictly positive, it is invertible.



## References

---

- D. Calvetti and E. Somersalo, *Mathematics of Data Science: A Computational Approach to Clustering and Classification*. SIAM, 2021, pp. 49-62.

# Mutidimensional scaling

## The basic ideas of MDS

---

Let  $X \in \mathbb{R}^{D \times n}$  be the high-dimensional data matrix

$$X = [x_1 \quad x_2 \quad \cdots \quad x_n],$$

and  $Y \in \mathbb{R}^{d \times n}$  be the dimension reduced data

$$Y = [y_1 \quad y_2 \quad \cdots \quad y_n].$$

MDS seeks to preserve that the dissimilarity (inner product) between all pairs of data points.

Set up the distance squared matrix

$$D^2 = [d_{ij}^2]_{n \times n},$$

where  $d_{ij} = \|x_i - x_j\|_2$ .

## Centering matrix

---

**Definition :** The **centering matrix** of size  $n$  is defined as

$$C_n = I_n - \frac{1}{n}J_n \in \mathbb{R}^{n \times n},$$

where  $I_n$  is the identity matrix of size  $n$  and  $J_n$  is an  $n \times n$  matrix of all 1's.

**Properties:**

- The eigenvalues of  $C_n$  are 0 of multiplicity 1 and 1 of multiplicity  $n - 1$ .
- $C_n$  is idempotent.

**Application:** Let  $X \in \mathbb{R}^{m \times m}$ .

- Each column of the product  $C_m X$  has a zero mean.
- Each row of the product  $X C_m$  has a zero mean.
- Each column and each row of double centering matrix  $C_m X C_m$  has a zero mean.

## Construct doubly centered squared distance matrix

Assume we centered the high-dimensional data such that

$$\sum_{i=1}^n x_i = \mathbf{0}.$$

The distance squared between two observations is

$$d_{ij}^2 = \|x_i - x_j\|_2^2 = \|x_i\|_2^2 + \|x_j\|_2^2 - 2x_i^\top x_j, \quad (1)$$

then we can write down the following formulas

$$\sum_{i=1}^n d_{ij}^2 = \text{Tr}(\mathbf{X}^\top \mathbf{X}) + n\|x_j\|_2^2 - 2\sum_{i=1}^n x_i^\top x_j = \text{Tr}(\mathbf{X}^\top \mathbf{X}) + n\|x_j\|_2^2 \quad (2)$$

$$\sum_{j=1}^n d_{ij}^2 = n\|x_i\|_2^2 + \text{Tr}(\mathbf{X}^\top \mathbf{X}) - 2x_i^\top \sum_{j=1}^n x_j = n\|x_i\|_2^2 + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \quad (3)$$

$$\sum_{i=1}^n \sum_{j=1}^n d_{ij}^2 = n \sum_{i=1}^n \|x_i\|_2^2 + n\text{Tr}(\mathbf{X}^\top \mathbf{X}) = 2n\text{Tr}(\mathbf{X}^\top \mathbf{X}). \quad (4)$$

## Construct doubly centered squared distance matrix (cont'd)

So combining the above three formulas  $((2) + (3))/n - (3)/n^2$ , we can get

$$\frac{1}{n} \left( \sum_{i=1}^n d_{ij}^2 + \sum_{i=1}^n d_{ij}^2 \right) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n d_{ij}^2 = \|x_i\|_2^2 + \|x_j\|_2^2. \quad (5)$$

We rewrite the (1)

$$\begin{aligned} x_i^\top x_j &= \frac{1}{2} (\|x_i\|_2^2 + \|x_j\|_2^2 - d_{ij}^2) \\ &= \frac{1}{2} \left( \frac{1}{n} \left( \sum_{i=1}^n d_{ij}^2 + \sum_{i=1}^n d_{ij}^2 \right) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n d_{ij}^2 - d_{ij}^2 \right) \\ &= -\frac{1}{2} \left( d_{ij} - \frac{1}{n} \left( \sum_{i=1}^n d_{ij}^2 + \sum_{i=1}^n d_{ij}^2 \right) + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n d_{ij}^2 \right). \end{aligned}$$

## Construct doubly centered squared distance matrix (cont'd)

---

We define the doubly centered distance matrix (inner product matrix)

$$K(D^2) = [k_{ij}]_{n \times n}$$

where

$$\begin{aligned} k_{ij} &= x_i^\top x_j \\ &= -\frac{1}{2} \left( d_{ij} - \frac{1}{n} \left( \sum_{i=1}^n d_{ij}^2 + \sum_{i=1}^n d_{ij}^2 \right) + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n d_{ij}^2 \right) \\ &= -\frac{1}{2} (C_n D^2 C_n)_{ij}. \end{aligned}$$

Hence,

$$K(D^2) = -\frac{1}{2} C_n D^2 C_n.$$

## Optimization problem

The optimization is

$$\min_{y_1, \dots, y_n} \sum_{ij} (y_i^\top y_j - k_{ij})^2 \implies \min_{\mathbf{Y}} \|\mathbf{Y}^\top \mathbf{Y} - \mathbf{K}(\mathbf{D}^2)\|_F^2.$$

Since  $\mathbf{K}(\mathbf{D}^2)$  is a symmetric matrix,

$$\begin{aligned} \mathbf{K}(\mathbf{D}^2) &= \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top \\ &= \begin{bmatrix} \mathbf{u}_1 & \cdots & \mathbf{u}_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_n^\top \end{bmatrix}, \end{aligned}$$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$ , and  $\mathbf{U}$  is orthonormal matrix.

So  $\mathbf{Y} = \mathbf{\Lambda}_d^{\frac{1}{2}} \mathbf{U}_d^\top$ , where  $\mathbf{\Lambda}_d$  is the diagonal matrix of first  $d$  eigenvalues and  $\mathbf{U}_d$  is the matrix of first  $d$  eigenvectors of  $\mathbf{K}(\mathbf{D}^2)$ .



## Classical MDS algorithm

---

- step 1: Identify  $k$  nearest neighbors (or  $\epsilon$  ball neighborhood) of each input data point and construct a neighborhood graph where edge lengths between points in a neighborhood are set as their Euclidean distances.
- step 2: **Compute Euclidean distances**,  $d_{ij}$ , that are  $\|x_i - x_j\|_2$ .
- step 3: Construct a matrix  $K(D^2) = -\frac{1}{2}CD^2C$ .
- step 4: Compute top  $d$  eigenvectors of  $K(D^2)$  which leads to the eigenvector matrix  $V \in \mathbb{R}^{n \times d}$  and the eigenvalue matrix  $\Lambda \in \mathbb{R}^{d \times d}$ .
- step 5: The coordinates of  $n$  points in the  $d$ -dimensional Euclidean space are given by the column vectors of  $Y = \Lambda^{1/2}V^\top$ .

## References

---

- I. Borg and P. J.F. Groenen, *Modern Multidimensional Scaling* (2ed ed.). Springer.
- K. Q. Weinberger and L. K. Saul (2006), Unsupervised Learning of image manifolds by semidefinite programming, *International Journal of Computer Vision*, 70(1), pp. 77-90. doi: 10.1007/s11263-005-4939-z.
- X. Zhen and Z. Liaangliang, The simplified expression of machine learning and multivariate statistical analysis based on the centering matrix, 2021.

## Isomap

## The basic ideas of Isomap

---

Isomap seeks to preserve the geodesic distances between all pairs of data points.

Define

$$D_G^2 = [d_G^2(i, j)]_{n \times n},$$

where  $d_G(i, j)$  is the geodesic distance between  $x_i$  and  $x_j$ .

Classical Isomap **assumes  $d_G$  has euclidean representation**, that is, there exists a set  $\{z_i\}_{i=1}^n \subset \mathbb{R}^D$  such that

$$\|z_i - z_j\|_2^2 = d_G^2(i, j).$$

Thus,  $D_G^2$  can be seen as distance squared matrix.

## Optimization problem

---

Then the algorithm is following the form of MDS:

$$K(D_G^2) = -\frac{1}{2}C_n D_G^2 C_n,$$

and the optimization problem is

$$\min_Y \|Y^\top Y - K(D_G^2)\|_F^2.$$

So

$$Y = \Lambda_d^{\frac{1}{2}} U_d^\top,$$

where  $\Lambda_d$  is the diagonal matrix of first  $d$  eigenvalues and  $U_d$  is the matrix of first  $d$  eigenvectors of  $K(D_G^2)$ .

## Additive constant problem

---

Let  $d$  be a measure of dissimilarity ( $d_{ii} = 0$ ;  $d_{ij} = d_{ji} \geq 0$ ) on a finite set  $I$  with  $n$  elements.

We consider the case where  $d^2$  has no euclidean representation, that is,  $W_{d^2}$  matrix is not positive semidefinite:

$$W_{d^2} = -\frac{1}{2}CD^2C,$$

with  $D^2$  is matrix with terms  $d_{ij}^2$ ,  $C$  is centering matrix

$$C = I - \frac{1}{n}jj^\top,$$

where  $I \in \mathbb{R}^{n \times n}$  is identity matrix,  $j \in \mathbb{R}^n$  is vector with all coordinates equal to 1.

Observe that

$$W_{d^2} = CW_{d^2} = W_{d^2}C, \quad \text{and that } W_{d^2}j = 0.$$

## Additive constant problem (cont'd)

---

We solve the "additive constant problem" which is find the smallest  $c^*$  such that the dissimilarity measure  $d_c$  defined by:

$$d_c(i, j) = \begin{cases} d_{ij} + c & , \text{ if } i \neq j \\ 0 & , \text{ if } i = j \end{cases}$$

has euclidean representation for all  $c \geq c^*$ .

Then define:

$$(\mathbf{D}_c^2)_{ij} := (d_{ij} + c)^2 = d_{ij}^2 + 2cd_{ij} + c^2 = (\mathbf{D}^2)_{ij} + 2c(\mathbf{D})_{ij} + c^2(\mathbf{J} - \mathbf{I})_{ij}.$$

The matrix  $\mathbf{W}$  associated to  $d_c$  define by:

$$\mathbf{W}_{d_c} = -\frac{1}{2}\mathbf{C}\mathbf{D}_c^2\mathbf{C} = \mathbf{W}_{d^2} + 2c\mathbf{W}_d + \frac{c^2}{2}\mathbf{C}, \quad (6)$$

where  $\mathbf{W}_d = -\frac{1}{2}\mathbf{C}\mathbf{D}\mathbf{C}$  with  $\mathbf{D}$  is matrix with terms  $d_{ij}$ .

## Additive constant problem (cont'd)

---

$d_c$  has an euclidean representation if:

$$\mathbf{x}^\top \mathbf{W}_{d_c} \mathbf{x} = \mathbf{x}^\top \mathbf{W}_{d^2} \mathbf{x} + 2c \mathbf{x}^\top \mathbf{W}_d \mathbf{x} + \frac{c^2}{2} \mathbf{x}^\top \mathbf{C} \mathbf{x}$$

is nonnegative for all  $\mathbf{x}$ .

### Theorem:

There exists a constant  $c^*$  such that the dissimilarity  $d_c$  has an euclidean representation for all  $c \geq c^*$ , and the euclidean representation of  $d_{c^*}$  defines a space with at most  $(n - 2)$  dimensions.



## Additive constant problem (cont'd)

---

*Proof:*

For a given  $x$ ,  $x^\top W_{d_c} x$  is function of  $c$  represented by a convex parabola; thus, to any  $x$  corresponds a number  $\alpha(x)$  such that:

$$x^\top W_{d_c} x \geq 0 \quad \text{if} \quad c \geq \alpha(x).$$

Observe that

$$\alpha(x) = \alpha(kx) \quad \text{for any number } k.$$

If the parabola is never negative,

$$\alpha(x) = -\infty.$$

For example,

$$\alpha(j) = -\infty.$$

## Additive constant problem (cont'd)

---

Because  $d^2$  has no euclidean representation, there is at least one  $x$  where  $x^\top W_{d^2} x < 0$  and for which  $\alpha(x)$  will be positive; thus, the number:

$$c^* = \sup_x \alpha(x) = \alpha(x^*)$$

is positive and such that:

$$x^\top W_{d_c} x \geq 0 \quad \text{for all } x \text{ and all } c \geq c^* \quad (7)$$

$$x^{*\top} W_{d_{c^*}} x^* = 0. \quad (8)$$

The eq.(7) proves that  $d_c$  has an euclidean representation for all  $c \geq c^*$ ; The eq.(8) shows that the representation of  $d_{c^*}$  defines a space with at most  $(n - 2)$  dimensions, since  $W_{d_{c^*}}$  has at least two independent eigenvectors ( $j$  and  $x^*$ ) associated to the eigenvalue zero.

## Calculation of $c^*$

---

As a consequence of eq.(7) and eq.(8),  $c^*$  and  $\mathbf{x}^*$  verify

$$\begin{aligned}W_{d_c^*} \mathbf{x}^* &= \mathbf{0} \\ \implies (W_{d^2} + 2c^* W_d + \frac{c^{*2}}{2} C) \mathbf{x}^* &= \mathbf{0} \\ \implies (W_{d^2} + 2c^* W_d + \frac{c^{*2}}{2} I) C \mathbf{x}^* &= \mathbf{0}. \quad (9)\end{aligned}$$

Let

$$2W_{d^2} \mathbf{x}^* = c^* \mathbf{y}, \quad (10)$$

we have the equation from eq.(9) since  $c^*$  is strictly positive :

$$\mathbf{y} + 4W_d \mathbf{x}^* + c^* C \mathbf{x}^* = \mathbf{0} \quad (11)$$

## Calculation of $c^*$ (cont'd)

---

Equalities eq.(10) and eq.(11) are gathered in the following:

$$\begin{pmatrix} \mathbf{0} & 2\mathbf{W}_{d^2} \\ -\mathbf{I} & -4\mathbf{W}_d \end{pmatrix} \begin{pmatrix} \mathbf{y} \\ \mathbf{C}\mathbf{x}^* \end{pmatrix} = c^* \begin{pmatrix} \mathbf{y} \\ \mathbf{C}\mathbf{x}^* \end{pmatrix},$$

which proves that  $c^*$  is an eigenvalue of the matrix

$$\mathbf{B} = \begin{pmatrix} \mathbf{0} & 2\mathbf{W}_{d^2} \\ -\mathbf{I} & -4\mathbf{W}_d \end{pmatrix}.$$

## Calculation of $c^*$ (cont'd)

Let  $a$  be an eigenvalue of  $B$  associated to the eigenvector  $\begin{pmatrix} z \\ t \end{pmatrix}$ :

$$\begin{cases} 2W_{d^2}t = az \\ -z - 4W_d t = at \end{cases},$$

we have

$$\begin{aligned} (W_{d^2} + 2aW_d + \frac{a^2}{2}I)t &= 0 \\ \implies (W_{d^2} + 2aW_d + \frac{a^2}{2}C)t &= 0 \\ \iff t^\top W_{d_a} t &= 0. \end{aligned}$$

Thus,  $\alpha(t) \geq a$ , which implies  $c^* \geq a$  because of the definition of  $c^*$ . The additive constant  $c^*$  is thus the largest eigenvalue of matrix  $B$ .

## Classical Isomap algorithm

---

- step 1: Identify  $k$  nearest neighbors (or  $\epsilon$  ball neighborhood) of each input data point and construct a neighborhood graph where edge lengths between points in a neighborhood are set as their Euclidean distances.
- step 2: Compute geodesic distances,  $d_G(i, j)$ , that are associated with the sum of edge weights along shortest paths between all pairs of points.
- step 3: Construct a matrix  $K(D_G^2) = -\frac{1}{2}CD_G^2C$ .
- step 4: Compute top  $d$  eigenvectors of  $K(D_G^2)$  which leads to the eigenvector matrix  $V \in \mathbb{R}^{n \times d}$  and the eigenvalue matrix  $\Lambda \in \mathbb{R}^{d \times d}$ .
- step 5: The coordinates of  $n$  points in the  $d$ -dimensional Euclidean space are given by the column vectors of  $Y = \Lambda^{1/2}V^\top$ .

## Kernel Isomap algorithm

---

- step 1: Identify  $k$  nearest neighbors (or  $\epsilon$  ball neighborhood) of each input data point and construct a neighborhood graph.
- step 2: Compute geodesic distances, that are associated with the sum of edge weights along shortest paths between all pairs of points.
- step 3: Construct a matrix  $K(D^2) = -\frac{1}{2}CD^2C$ .
- step 4: Compute the **largest eigenvalue**  $c^*$  of the matrix

$$\begin{bmatrix} \mathbf{0} & 2K(D^2) \\ -I & -4K(D) \end{bmatrix} \quad (12)$$

and construct a **Mercer kernel matrix**  $\tilde{K} = \tilde{K}(\tilde{D}^2)$ , that is,

$$\tilde{K} = K(D^2) + 2cK(D) + \frac{1}{2}c^2C,$$

where  $\tilde{K}$  is guaranteed to be positive semidefinite for  $c \geq c^*$ .

## Kernel Isomap algorithm (cont'd)

---

- step 5: Compute top  $d$  eigenvectors of  $\tilde{K}$  which leads to the eigenvector matrix  $V \in \mathbb{R}^{n \times d}$  and the eigenvalue matrix  $\Lambda \in \mathbb{R}^{d \times d}$ .
- step 6: The coordinates of  $n$  points in the  $d$ -dimensional Euclidean space are given by the column vectors of  $Y = \Lambda^{1/2} V^\top$ .



## References

---

- I. Borg and P. J.F. Groenen. *Modern Mutidimensional Scaling* (2ed ed.). Springer.
- K. Q. Weinberger and L. K. Saul (2006), Unsupervised Learning of image manifolds by semidefinite programming, *International Journal of Computer Vision*, 70(1), pp. 77-90. doi: 10.1007/s11263-005-4939-z.
- X. Zhen and Z. Liaangliang. The simplified expression of machine learning and multivariate statistical analysis based on the centering matrix, 2021.
- C. Heeyoul and C. Seungjin. Robust kernel Isomap, *ScienceDirect*, 2006, pp. 853-862.
- F. Cailliez (1983). The analytical solution of the additive constant problem, *PSYCHOMETRIKA*, 40(2), pp. 305-308.
- A. Saxena, A. Gupta, and A. Mukerjee (2004). Non-linear dimensionality reduction by locally linear Isomaps. *Lecture Notes in Computer Science*(pp. 1038-1043). Research Gate.

## References

---

- J. B. Tenenbaum, V. D. Silva, and J. C. Langford (2000). A global framework for nonlinear dimensionality reduction. *SCIENCE*, 290, pp. 2319-2322.

## Diffusion maps

## The basic ideas of Diffusion maps

---

Let  $X \in \mathbb{R}^{D \times n}$  be the high-dimensional data matrix

$$X = [x_1 \quad x_2 \quad \cdots \quad x_n],$$

and  $Y \in \mathbb{R}^{d \times n}$  be the dimension reduced data

$$Y = [y_1 \quad y_2 \quad \cdots \quad y_n].$$

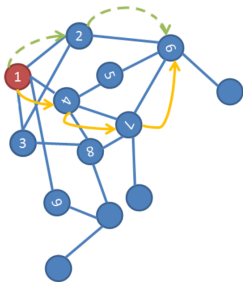
Diffusion map is construct the reduction data from intrinsic local geometric relationship, i.e,

$$\|x_i - x_j\|_D \text{ is small} \implies \|y_i - y_j\|_2 \text{ is small},$$

where  $\|x_i - x_j\|_D$  is diffusion distance between  $x_i$  and  $x_j$ .

## Connectivity

Suppose we take a random walk on our data, jumping between data points in feature space.



**Figure:** A random walk on a data set. Each "jump" has a probability associated with it. The dashed path between nodes 1 and 6 requires two jumps (i.e., two time units) with the probability along the path being  $p(\text{node 1, node 2}) p(\text{node 2, node 6})$ .

## Connectivity (cont'd)

---

The connectivity between two data points,  $x$  and  $y$ , is defined as the probability of jumping from  $x$  to  $y$  in one step of the random walk, and is

$$\text{connectivity}(x, y) = p(x, y). \quad (13)$$

It is useful to express this connectivity in terms of a non-normalized likelihood function,  $k$ , known as the diffusion kernel:

$$\text{connectivity}(x, y) \propto k(x, y). \quad (14)$$

## Connectivity (cont'd)

---

The diffusion kernel satisfies the following properties:

1.  $k$  is symmetric :  $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{y}, \mathbf{x})$
2.  $k$  is positivity preserving :  $k(\mathbf{x}, \mathbf{y}) \geq 0$

The latter property is specific to the diffusion kernel and allows it to be interpreted as a scaled probability, so that

$$\frac{1}{d_X} \sum_{\mathbf{y} \in X} k(\mathbf{x}, \mathbf{y}) = 1.$$

The relation between the kernel function and the connectivity is then

$$\text{connectivity}(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}, \mathbf{y}) = \frac{1}{d_X} k(\mathbf{x}, \mathbf{y})$$

with  $\frac{1}{d_X}$  the normalization constant.

## Connectivity (cont'd)

---

Define a row-normalized diffusion matrix,  $\mathbf{P}$ , with entries

$$P_{ij} = p(x_i, x_j).$$

By taking powers of the diffusion matrix, we can increase the number of steps taken. For example, take a  $2 \times 2$  diffusion matrix,

$$\mathbf{P} = \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix}.$$

Each element,  $P_{ij}$ , is the probability of jumping between data points  $i$  and  $j$ . When  $\mathbf{P}$  is squared, it becomes

$$\mathbf{P}^2 = \begin{bmatrix} p_{11}p_{11} + p_{12}p_{21} & p_{12}p_{22} + p_{11}p_{12} \\ p_{21}p_{11} + p_{22}p_{21} & p_{22}p_{22} + p_{21}p_{12} \end{bmatrix}.$$

When making two jumps, these  $P_{ij}$  are all the paths from point  $i$  to point  $j$ . Similarly,  $\mathbf{P}^t$  sum all paths of length  $t$  from point  $i$  to point  $j$ .



## Diffusion Process

---

With increased values of  $t$  (i.e. as the diffusion process), the probability of following a path along the underlying geometric structure of the data set increases.



**Figure:** Paths along the true geometric structure of the data set have high probability.

## Diffusion Distance

---

The diffusion distance between  $x_i$  and  $x_j$  is defined by

$$\begin{aligned} D_t(x_i, x_j)^2 &= \sum_{u \in X} |p_t(x_i, u) - p_t(x_j, u)|^2 \\ &= \sum_k |P_{ik}^t - P_{jk}^t|^2 \end{aligned}$$

The diffusion distance is small if there are many high probability paths of length  $t$  between two points.

## Diffusion map

---

The diffusion distance in data space simply becomes the Euclidean distance in this new diffusion space. With this in mind, we define and examine the mapping

$$\mathbf{y}_i := \begin{bmatrix} p_t(\mathbf{x}_i, \mathbf{x}_1) \\ p_t(\mathbf{x}_i, \mathbf{x}_2) \\ \vdots \\ p_t(\mathbf{x}_i, \mathbf{x}_n) \end{bmatrix} = \mathbf{P}_{i*}^\top. \quad (15)$$

For this map, the Euclidean distance between  $\mathbf{y}_i$  and  $\mathbf{y}_j$  is

$$\begin{aligned} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 &= \sum_{\mathbf{u} \in X} |p_t(\mathbf{x}_i, \mathbf{u}) - p_t(\mathbf{x}_j, \mathbf{u})|^2 \\ &= \sum_k |\mathbf{P}_{ik}^t - \mathbf{P}_{jk}^t|^2 = D_t(\mathbf{x}_i, \mathbf{x}_j)^2. \end{aligned}$$

Note that the dimension of the mapped data is still the sample size  $n$ .

## Diffusion map (cont'd)

---

**Lemma 1:** Suppose  $K \in \mathbb{R}^{n \times n}$  is a kernel matrix such that  $K[i, j] = k(i, j)$ . A diagonal matrix,  $D$ , normalizes the rows of  $K$  to produce a diffusion matrix

$$P = D^{-1}K. \quad (16)$$

Then, matrix  $P'$ , defined as

$$P' = D^{1/2}PD^{-1/2}, \quad (17)$$

1. is symmetric.
2. has the same eigenvalues as  $P$ .
3. the eigenvectors,  $x'_k$  of  $P'$  are multiplied by  $D^{-1/2}$  and  $D^{1/2}$  to get the left ( $u^\top P = \lambda u^\top$ ) and right eigenvectors ( $Pv = \lambda v$ ) of  $P$  respectively.

## Diffusion map (cont'd)

---

**proof :** Substitute (16) into (17) to obtain

$$\mathbf{P}' = \mathbf{D}^{-1/2} \mathbf{K} \mathbf{D}^{-1/2}. \quad (18)$$

As  $\mathbf{K}$  is symmetric,  $\mathbf{P}'$  will also be symmetric.

We can make  $\mathbf{P}$  the subject of the equation in (17),

$$\mathbf{P} = \mathbf{D}^{-1/2} \mathbf{P}' \mathbf{D}^{1/2}. \quad (19)$$

As  $\mathbf{P}'$  is symmetric, by spectrum analysis such that

$$\mathbf{P}' = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^{\top} \quad (20)$$

where  $\mathbf{\Lambda}$  is a diagonal matrix containing the eigenvalues of  $\mathbf{P}'$  and  $\mathbf{S}$  is a matrix with the orthonormal eigenvectors of  $\mathbf{P}'$  as columns.

## Diffusion map (cont'd)

---

Substituting (20) into (19),

$$\mathbf{P} = \mathbf{D}^{-1/2} \mathbf{S} \mathbf{\Lambda} \mathbf{S}^\top \mathbf{D}^{1/2} \quad (21)$$

$$= \mathbf{D}^{-1/2} \mathbf{S} \mathbf{\Lambda} \mathbf{S}^{-1} \mathbf{D}^{1/2} \quad (22)$$

$$= (\mathbf{D}^{-1/2} \mathbf{S}) \mathbf{\Lambda} (\mathbf{S}^{-1} \mathbf{D}^{1/2}) \quad (23)$$

$$= \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}, \quad (24)$$

where  $\mathbf{Q} = \mathbf{D}^{-1/2} \mathbf{S}$ . Therefore, the eigenvalues of  $\mathbf{P}$  and  $\mathbf{P}'$  are the same.

1. The right eigenvectors of  $\mathbf{P}$  are the columns of  $\mathbf{Q}$
2. The left eigenvectors of  $\mathbf{P}$  are the rows of  $\mathbf{Q}^{-1}$

## Diffusion map (cont'd)

---

Observe that the eigenvectors of  $P$  can be given in terms of the eigenvectors  $x'_k$  of  $P'$ . The right eigenvectors of  $P$  are

$$v_k = D^{-1/2} x'_k \quad (25)$$

and the left eigenvectors are

$$u_k = D^{1/2} x'_k. \quad (26)$$

by SVD, we obtain

$$P = \sum_k \lambda_k v_k u_k^\top. \quad (27)$$

Eq.(27) expresses each row of the diffusion matrix in terms of a new basis:  $u_k$ , the left eigenvectors of the diffusion matrix.

## Diffusion map (cont'd)

---

In this new coordinate system in  $\mathbb{R}^n$ , a row  $i$  of  $P$ , which is  $(y_i)$  is represented by the point

$$M_i = \begin{bmatrix} \lambda_1 v_1[i] \\ \lambda_2 v_2[i] \\ \vdots \\ \lambda_n v_n[i] \end{bmatrix},$$

where  $v_n[i]$  is the  $i$ -th component of the  $n$ -th right eigenvector. Map to the  $d$ -dimensional diffusion space at time  $t$ , using the  $d$  dominant eigenvectors and -values as

$$y_i = \begin{bmatrix} \lambda_1^t v_1[i] \\ \lambda_2^t v_2[i] \\ \vdots \\ \lambda_d^t v_d[i] \end{bmatrix}.$$



## References

---

- J. Porte, B. M. Herbst, W. Hereman, and S. J. Walt. An Introduction to diffusion maps, 2008.

## Laplacian Eigenmaps (LE)

## The basic ideas of LE

---

Let  $\mathbf{X}$  be a dataset

$$\mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_n] \in \mathbb{R}^{D \times n},$$

and  $\mathbf{Y}$  be the dimension reduced data

$$\mathbf{Y} = [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \cdots \quad \mathbf{y}_n] \in \mathbb{R}^{d \times n}.$$

LE construct the reduction data from intrinsic local geometric relationship, i.e,

$$\|\mathbf{x}_i - \mathbf{x}_j\|_2 \text{ is small} \implies \|\mathbf{y}_i - \mathbf{y}_j\|_2 \text{ is small},$$

where  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are original high dimensional data, and  $\mathbf{y}_i$  and  $\mathbf{y}_j$  are lower dimensional data.

## Constructing the adjacency graph

---

- 1)  $\epsilon$ -neighborhoods (parameter  $\epsilon \in \mathbb{R}$ ). Nodes  $i$  and  $j$  are connected by an edge if  $\|x_i - x_j\|^2 < \epsilon$  where the norm is the usual Euclidean norm.
- 2)  $k$  nearest neighbors (parameter  $k \in \mathbb{N}$ ). Node  $i$  and  $j$  are connected by an edge if  $i$  is among  $k$  nearest neighbors of  $j$  or  $j$  is among  $k$  nearest neighbors of  $i$ . Note that this relation is symmetric.
- 3) Fully connected.

## Choosing the weights

---

Let  $W \in \mathbb{R}^{n \times n}$ , we have two variations for weighting the edges:

- 1) Heat kernel (parameter  $t \in \mathbb{R}$ ). If nodes  $i$  and  $j$  are connected, put

$$W_{ij} = e^{-\frac{\|x_i - x_j\|^2}{t}};$$

otherwise, put  $W_{ij} = 0$ .

- 2) Simple-minded (no parameters ( $t = \infty$ )).

$$\begin{cases} W_{ij} = 1, & \text{if vertices } i \text{ and } j \text{ are connected by an edge.} \\ W_{ij} = 0, & \text{if vertices } i \text{ and } j \text{ are not connected by an edge.} \end{cases}$$

This simplification avoids the need to choose  $t$ .

## Optimization problem

This motivates us define the cost function  $\Phi(\mathbf{Y})$ :

$$\begin{aligned}\Phi(\mathbf{Y}) &= \sum_{i,j} \mathbf{W}_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \\&= \sum_{i,j} \mathbf{W}_{ij} (\|\mathbf{y}_i\|^2 + \|\mathbf{y}_j\|^2 - 2\mathbf{y}_i^\top \mathbf{y}_j) \\&= \sum_i \|\mathbf{y}_i\|^2 \sum_j \mathbf{W}_{ij} + \sum_j \|\mathbf{y}_j\|^2 \sum_i \mathbf{W}_{ij} - 2 \sum_{ij} \mathbf{W}_{ij} \mathbf{y}_i^\top \mathbf{y}_j \\&= 2 \sum_i \|\mathbf{y}_i\|^2 \sum_j \mathbf{W}_{ij} - 2 \sum_{ij} \mathbf{W}_{ij} \mathbf{y}_i^\top \mathbf{y}_j \\&= 2 \sum_i \|\mathbf{y}_i\|^2 d_{ii} - 2 \sum_{ij} \mathbf{W}_{ij} \mathbf{y}_i^\top \mathbf{y}_j,\end{aligned}$$

where  $d_{ii} = \sum_{j=1}^n \mathbf{W}_{ij}$ .

## Optimization problem

---

$$\begin{aligned}\Phi(\mathbf{Y}) &= 2\text{tr}(\mathbf{Y}\mathbf{D}\mathbf{Y}^\top) - 2\text{tr}(\mathbf{Y}\mathbf{W}\mathbf{Y}^\top) \\ &= 2\text{tr}(\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top) \\ &= 2\text{tr}(\mathbf{Y}\mathcal{L}\mathbf{Y}^\top),\end{aligned}$$

where  $\mathbf{D} = \text{diag}(d_{ii}) \in \mathbb{R}^{n \times n}$ ,  $\mathcal{L} = \mathbf{D} - \mathbf{W} \in \mathbb{R}^{n \times n}$  is called Laplacian matrix.

So the optimization problem is

$$\min_{\mathbf{Y} \in \mathbb{R}^{d \times n}} \text{tr}(\mathbf{Y}\mathcal{L}\mathbf{Y}^\top) \quad \text{subject to } \mathbf{Y}\mathbf{D}\mathbf{1}_n = 0, \mathbf{Y}\mathbf{D}\mathbf{Y}^\top = \mathbf{I}_{d \times d},$$

where  $\mathbf{1}_n \in \mathbb{R}^n$  is a vector with all 1's and  $\mathbf{I}$  is identity matrix.

## Optimization problem

---

Define lagrange function

$$L_{\lambda, \Lambda}(\mathbf{Y}) = \text{tr}(\mathbf{Y} \mathcal{L} \mathbf{Y}^\top) + \lambda^\top \mathbf{Y} \mathbf{D} \mathbf{1}_n + \text{tr}(\Lambda(\mathbf{I} - \mathbf{Y} \mathbf{D} \mathbf{Y}^\top)),$$

where  $\lambda \in \mathbb{R}^d$  and  $\Lambda \in \mathbb{R}^{d \times d}$ , diagonal matrix, are lagrange multipliers.

Find

$$\frac{\partial L}{\partial \mathbf{Y}^\top} = 2\mathcal{L} \mathbf{Y}^\top + \mathbf{D} \mathbf{1}_n \lambda^\top - 2\mathbf{D} \mathbf{Y}^\top \Lambda = 0,$$

we obtained

$$\begin{aligned} \lambda &= \mathbf{0} \\ \mathcal{L} \mathbf{Y}^\top &= \mathbf{D} \mathbf{Y}^\top \Lambda. \end{aligned}$$

Hence

$$\mathbf{D}^{-1} \mathcal{L} \mathbf{Y}^\top = \mathbf{Y}^\top \Lambda.$$



## Eigenvalue problem

From above results, we get an eigenvalue problem of  $D^{-1}\mathcal{L}$ .  
Observing that

$$\begin{aligned}\mathcal{L}\mathbf{1}_N &= (D - W)\mathbf{1}_n \\ &= D\mathbf{1}_n - W\mathbf{1}_n \\ &= \begin{bmatrix} d_{11} \\ d_{22} \\ \vdots \\ d_{nn} \end{bmatrix} - \begin{bmatrix} \sum_j W_{1j} \\ \sum_j W_{2j} \\ \vdots \\ \sum_j W_{nj} \end{bmatrix} \\ &= \mathbf{0}_n,\end{aligned}$$

the  $\mathbf{1}_n$  is eigenvector corresponding to 0 eigenvalue conflicts with first constraint, so we select the eigenvector corresponding to the second smallest eigenvalue to the  $(d + 1)$ -th smallest eigenvalue.

## References

---

- M. Belkin and P. Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering, 2001.

### Locally linear embedding (LLE)

## The basic ideas of LLE

---

Let  $X$  be a dataset

$$X = [x_1 \quad x_2 \quad \cdots \quad x_n] \in \mathbb{R}^{D \times n},$$

and  $Y$  be the dimension reduced data

$$Y = [y_1 \quad y_2 \quad \cdots \quad y_n] \in \mathbb{R}^{d \times n}.$$

LLE construct the dimension reduction data from intrinsic local linear relationship. For example,

$$x_1 = c_1 x_2 + c_2 x_3 + c_3 x_4,$$

where  $c_1$ ,  $c_2$ , and  $c_3$  are called weight coefficients. After using LLE to reduce dimension,

$$y_1 \approx c_1 y_2 + c_2 y_3 + c_3 y_4.$$

We choose  $k \geq d$ , and use  $K$ -nearest neighbor(KNN) algorithm to find  $k$ -nearest neighbor of each  $x_j$ .

Denoting the neighbor matrices of  $x_j$  by

$$N_j = knn(x_j, k) = [x_{j_1} \quad \cdots \quad x_{j_k}] \in \mathbb{R}^{D \times k}, \quad 1 \leq j \leq n,$$

where  $j_i \in \{1, 2, \dots, n\}$ .

Then we also define coefficient matrix  $W \in \mathbb{R}^{k \times n}$

$$W = [w_1 \quad w_2 \quad \cdots \quad w_n],$$

where weights  $w_j \in \mathbb{R}^k$  is coefficient vector, summarizing the contribution of the  $x_{j_i}$  to the  $j$ th reconstruction. That is,

$$x_j = N_j w_j = [x_{j_1} \quad x_{j_2} \quad \cdots \quad x_{j_k}] w_j, \quad 1 \leq j \leq n.$$

## LLE algorithm

---

**Step 1 .** Find the coefficient matrix, i.e, set  $w_j$  to be a solution of the problems

$$\min_{w_j \in \mathbb{R}^k} \sum_{j=1}^n \|x_j - \sum_{i=1}^k W_{ij} x_{j_i}\|_2^2 \quad \text{subject to} \quad \sum_{i=1}^k W_{ij} = 1, \quad 1 \leq j \leq n. \quad (P1)$$

**Step 2 .** With  $w_j$  given from Step 1, set  $Y$  to be a solution of the problem

$$\min_{\{y_i\}_{i=1}^n \subset \mathbb{R}^d} \sum_{j=1}^n \|y_j - \sum_{i=1}^k W_{ij} y_{j_i}\|_2^2 \quad \text{subject to} \quad \sum_{j=1}^n y_j = \mathbf{0} \text{ and } YY^\top = n \times I_d, \quad (P2)$$

where  $I_d$  is  $d \times d$  the identity matrix.

## Optimizing the weights

---

Start by writing the latter using matrix notation,

$$\begin{aligned}\sum_{j=1}^n \left\| \mathbf{x}_j - \sum_{i=1}^k \mathbf{w}_{ij} \mathbf{x}_{j_i} \right\|^2 &= \sum_{j=1}^n \left\| \mathbf{1} \times \mathbf{x}_j - \sum_{i=1}^k \mathbf{w}_{ij} \mathbf{x}_{j_i} \right\|^2 \\ &= \sum_{j=1}^n \left\| \sum_{i=1}^k \mathbf{w}_{ij} \mathbf{x}_j - \sum_{i=1}^k \mathbf{w}_{ij} \mathbf{x}_{j_i} \right\|^2 \\ &= \sum_{j=1}^n \left\| \sum_{i=1}^k \mathbf{w}_{ij} (\mathbf{x}_j - \mathbf{x}_{j_i}) \right\|^2. \\ &= \sum_{j=1}^n \left\| (\tilde{\mathbf{X}}_j - \mathbf{N}_j) \mathbf{w}_j \right\|^2,\end{aligned}$$

where  $\tilde{\mathbf{X}}_j = [\mathbf{x}_j \quad \cdots \quad \mathbf{x}_j] \in \mathbb{R}^{D \times k}$ .

## Optimizing the weights

---

It follows that

$$\begin{aligned}\sum_{j=1}^n \|(\tilde{\mathbf{X}}_j - \mathbf{N}_j)\mathbf{w}_j\|^2 &= \sum_{j=1}^n \mathbf{w}_j^\top (\tilde{\mathbf{X}}_j - \mathbf{N}_j)^\top (\tilde{\mathbf{X}}_j - \mathbf{N}_j) \mathbf{w}_j \\ &= \sum_{j=1}^n \mathbf{w}_j^\top \mathbf{S}_j \mathbf{w}_j,\end{aligned}\tag{28}$$

where  $\mathbf{S}_j = (\tilde{\mathbf{X}}_j - \mathbf{N}_j)^\top (\tilde{\mathbf{X}}_j - \mathbf{N}_j) \in \mathbb{R}^{k \times k}$  are *local covariance matrix*.  
And the constraint

$$\sum_{i=1}^k \mathbf{W}_{ij} = \mathbf{1}_k^\top \mathbf{w}_j = 1,$$

where  $\mathbf{1}_k \in \mathbb{R}^k$  is vector with all components are 1.



## Optimizing the weights

---

Hence, the problems (Q1) can be restated componentwise as

$$\min_{\mathbf{w}_j} \mathbf{w}_j^\top \mathbf{S}_j \mathbf{w}_j \quad \text{subject to} \quad \mathbf{1}_k^\top \mathbf{w}_j = 1, \quad 1 \leq j \leq n.$$

Define a Lagrange function

$$L_\lambda(\mathbf{w}_j) = \mathbf{w}_j^\top \mathbf{S}_j \mathbf{w}_j + \lambda(\mathbf{1}_k^\top \mathbf{w}_j - 1),$$

where  $\lambda \in \mathbb{R}$  is a Lagrange multiplier.

And look for its critical points, i.e, vector  $\mathbf{w}_j$  such that

$$\nabla_{\mathbf{w}_j} L_\lambda(\mathbf{w}_j) = 2\mathbf{S}_j \mathbf{w}_j + \lambda \mathbf{1}_k = 0.$$

We have

$$\mathbf{w}_j = -\frac{\lambda \mathbf{S}_j^{-1} \mathbf{1}_k}{2}. \quad (29)$$

## Optimizing the weights

---

Then focusing on constraint,

$$\mathbf{1}_k^\top \mathbf{w}_j = -\frac{\lambda \mathbf{1}_k^\top \mathbf{S}_j^{-1} \mathbf{1}_k}{2} = 1,$$

that gives

$$\lambda = -\frac{2}{\mathbf{1}_k^\top \mathbf{S}_j^{-1} \mathbf{1}_k}.$$

Hence, substituting the expression of  $\lambda$  in terms of  $\mathbf{w}_j$  in (29),

$$\mathbf{w}_j = \frac{\mathbf{S}_j^{-1} \mathbf{1}_k}{\mathbf{1}_k^\top \mathbf{S}_j^{-1} \mathbf{1}_k},$$

where the denominator is sum of entries of  $\mathbf{S}_j^{-1}$ , and numerator is sum of column vectors of  $\mathbf{S}_j^{-1}$ . It has a unique solution if the *local covariance matrix*  $\mathbf{S}_j$  is invertible.

## Optimizing the weights

---

On the other hand, if  $S_j$  is singular, it may have multiple solutions. In this situation it is suggested that we add a small regularization term  $\epsilon I$  to  $S_j$ , and set

$$w_j = \frac{(S_j + \epsilon I)^{-1} \mathbf{1}_k}{\mathbf{1}_k^\top (S_j + \epsilon I)^{-1} \mathbf{1}_k}.$$

## Mapping to lower dimensional $\mathbb{R}^d$

We begin by introducing a sparse matrix  $\widetilde{\mathbf{W}} \in \mathbb{R}^{n \times n}$ ,

$$\widetilde{\mathbf{W}}_{j,i} = \begin{cases} \mathbf{W}_{ij} & , \text{ if } \mathbf{x}_{j_i} \text{ is column of } N_j \\ \mathbf{0} & , \text{ if } \mathbf{x}_{j_i} \text{ is not column of } N_j \end{cases} \quad , \quad \text{for } 1 \leq i \leq k \text{ and } 1 \leq j \leq n.$$

That is  $\widetilde{\mathbf{W}}_{ij} = 0$  if  $x_i$  does not belong to the set of neighbors of  $x_j$ .  
Observing that

$$\sum_{i=1}^k \mathbf{W}_{ij} \mathbf{y}_{j_i} = \mathbf{Y} \widetilde{\mathbf{W}}_j,$$

so we can rewrite the cost function,

$$\sum_{j=1}^N \left\| \mathbf{y}_j - \sum_{i=1}^k \mathbf{W}_{ij} \mathbf{y}_{j_i} \right\|^2 = \sum_{j=1}^N \left\| \mathbf{Y}(\mathbf{I}_j - \widetilde{\mathbf{W}}_j) \right\|^2,$$

where  $\mathbf{I}_j$  is  $j$ th column of identity matrix  $\mathbf{I} \in \mathbb{R}^{d \times d}$ .

## Mapping to lower dimensional $\mathbb{R}^d$

---

It follows that,

$$\begin{aligned}\sum_{j=1}^N \|\mathbf{Y}(\mathbf{I}_j - \widetilde{\mathbf{W}}_j)\|^2 &= \text{tr}(\mathbf{Y}(\mathbf{I} - \widetilde{\mathbf{W}})(\mathbf{I} - \widetilde{\mathbf{W}})^\top \mathbf{Y}^\top) \\ &= \text{tr}(\mathbf{Y}\mathbf{M}\mathbf{Y}^\top)\end{aligned}$$

where

$$\mathbf{M} = (\mathbf{I} - \widetilde{\mathbf{W}})(\mathbf{I} - \widetilde{\mathbf{W}})^\top \in \mathbb{R}^{n \times n}.$$

Hence, the optimization problem (Q2) can be restated as,

$$\min_{\mathbf{Y}} \text{tr}(\mathbf{Y}\mathbf{M}\mathbf{Y}^\top) \quad \text{subject to } \mathbf{Y}\mathbf{Y}^\top = N \times \mathbf{I} \text{ and } \sum_{j=1}^n y_j = \mathbf{0}.$$

## Mapping to lower dimensional $\mathbb{R}^d$

---

We define a Lagrange function

$$L_{\lambda, \Lambda}(\mathbf{Y}) = \text{tr}(\mathbf{Y}\mathbf{M}\mathbf{Y}^\top) - \text{tr}(\Lambda(\mathbf{Y}\mathbf{Y}^\top - n\mathbf{I})) - \lambda^\top \mathbf{Y}\mathbf{1}_n,$$

where  $\Lambda \in \mathbb{R}^{d \times d}$  is a diagonal matrix and  $\lambda \in \mathbb{R}^d$ . And look for its critical points, i.e, matrix  $\mathbf{Y}$  such that

$$\frac{\partial L}{\partial \mathbf{Y}} = 2\mathbf{M}\mathbf{Y}^\top - 2\mathbf{Y}^\top \Lambda - \mathbf{1}_n \lambda^\top = 0.$$

We get,

$$\lambda = \mathbf{0} \quad \text{and} \quad \mathbf{M}\mathbf{Y}^\top = \mathbf{Y}^\top \Lambda.$$

It can be seen that  $\mathbf{Y}^\top$  is actually a matrix composed of eigenvectors of  $\mathbf{M}$ . In order to reduce the data to  $d$  dimension, **we only need to take the eigenvector corresponding to the second smallest eigenvalue to the  $(d + 1)$ -th smallest eigenvalue.**

## References

---

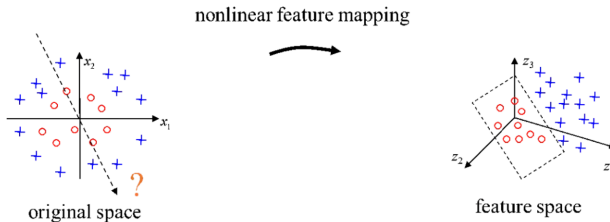
- S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding, *SCIENCE*, 290, pp. 2323-2325.
- S. T. Roweis and L. K. Saul. An introduction to locally linear embedding, 2001.
- L. Liren. Avoiding unwanted results in locally linear embedding: A new understanding of regularization, 2021.

## Kernel method



## Kernel method (Kernel trick)

The role of Kernel trick in machine learning is to hope that the data after an nonlinear projection can be more separated in a higher-dimensional space when different types of data cannot be separated by linear classifiers in the original space.



## Kernel function

---

**Definition :** For all  $x_i$  and  $x_j$  in the input space  $\mathcal{X}$ , certain functions  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  can be expressed as an inner product in other space  $\mathcal{V}$ . The function  $k$  is often referred to as a *kernel* or a *kernel function*. That is,

$$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle,$$

where  $\phi : \mathcal{X} \rightarrow \mathcal{V}$  is a nonlinear mapping.

**Mercer's condition (Discrete analog) :** A matrix  $K \in \mathbb{R}^{n \times n}$ ,  $K_{ij} = k(x_i, x_j)$ , which satisfies for all vectors  $g \in \mathbb{R}^n$ , the property,

$$\langle g, Kg \rangle = g^\top K g \geq 0.$$

**Mercer's theorem :** An implicitly defined function  $\phi$  exists whenever the space  $\mathcal{X}$  can be equipped with a suitable measure ensuring the kernel function  $k$  satisfies Mercer's condition.

## Commonly used kernel functions

---

- ① Linear kernel :

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle .$$

- ② Polynomial kernel :

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d ,$$

where  $d \in \mathbb{Z}^+$ .

- ③ Gaussian Radial Basis Function kernel (RBF) :

$$k(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}} ,$$

where  $\sigma \in \mathbb{R} - \{0\}$ .

## Kernel-based PCA

---

Let  $\Phi$  be a data matrix in the feature space,

$$\Phi = [\phi(x_1) \quad \phi(x_2) \quad \cdots \quad \phi(x_n)],$$

where  $\phi(\cdot)$  is a nonlinear mapping onto feature space.

Define the covariance matrix

$$C = \frac{1}{n}(\Phi H)(\Phi H)^\top, \quad (30)$$

where

$$H = I_n - \frac{1}{n}J_n \in \mathbb{R}^{n \times n},$$

is centering matrix,  $I_n$  is identity matrix of size  $n$  and  $J_n$  is an  $n \times n$  matrix of all 1's.

## Eigenvectors of $\tilde{C}$

---

Define the Mercer kernel matrix (using linear kernel here)

$$K = (\Phi H)^\top (\Phi H) = H \Phi^\top \Phi H \in \mathbb{R}^{n \times n}$$

The eigendecomposition of  $K$  is of the form

$$K = V \Lambda V^\top, \quad (31)$$

where  $V = [v_1 \ \cdots \ v_n]$  is orthonormal and  $\Lambda = \text{diag}(\lambda_i)$ ,  $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$ .

Pre-multiplying both sides of Eq.(31) by  $\Phi H$  leads to

$$(\Phi H) K V = (\Phi H) V \Lambda. \quad (32)$$

Then, observing  $H$  is symmetric, we have

$$\Phi H H \Phi^\top \Phi H V = \Phi H V \Lambda. \quad (33)$$

## Eigenvectors of $\tilde{C}$ (cont'd)

Thus, Eq.(33) can be written as

$$\Phi H H \Phi^\top \Phi H V = n C \Phi H V = \Phi H V \Lambda \quad (34)$$

$$n C U = U \Lambda, \quad (35)$$

where  $U = \Phi H V$  is the eigenvector matrix of  $C$ . That is,

$$U = \Phi H \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \quad (36)$$

$$= \begin{bmatrix} \Phi H v_1 & \cdots & \Phi H v_n \end{bmatrix} \quad (37)$$

$$= \begin{bmatrix} u_1 & \cdots & u_n \end{bmatrix} \quad (38)$$

where  $u_i = \Phi H v_i$ .

Hence,  $i$ -th best projection direction (maximizing the variance) in feature space,  $u_i$ , the eigenvector of  $nC$ , is found.

## Coordinates in the feature space

---

Normalized the  $u_i$

$$\begin{aligned}\tilde{u}_i &= \frac{1}{\|\Phi H v_i\|} \Phi H v_i = \frac{1}{\sqrt{v_i^\top H \Phi^\top \Phi H v_i}} \Phi H v_i \\ &= \frac{1}{\sqrt{v_i^\top H \Phi^\top \Phi H v_i}} \Phi H v_i = \frac{1}{\sqrt{v_i^\top K v_i}} \Phi H v_i \\ &= \frac{1}{\sqrt{v_i^\top \lambda_i v_i}} \Phi H v_i = \frac{1}{\sqrt{\lambda_i}} \Phi H v_i,\end{aligned}$$

where  $\lambda_i$  is the corresponding value of  $v_i$ .

## Coordinates in the feature space (cont'd)

---

Centering  $\Phi$  and projecting onto the unit vector  $\tilde{u}_i$  leads to the eigenvector  $v_i$  scaled by  $\lambda^{1/2}$ , i.e.,

$$\begin{aligned} \langle \Phi H, \tilde{u}_i \rangle &= (\Phi H)^\top \tilde{u}_i \\ &= (\Phi H)^\top \frac{1}{\sqrt{\lambda_i}} \Phi H v_i \\ &= \frac{1}{\sqrt{\lambda_i}} \mathbf{H} \Phi^\top \Phi H v_i \\ &= \frac{1}{\sqrt{\lambda_i}} \mathbf{K} v_i \\ &= \frac{1}{\sqrt{\lambda_i}} \lambda_i v_i \\ &= \lambda_i^{1/2} v_i. \end{aligned}$$



## References

---

- X. Zhen and Z. Liaangliang. The simplified expression of machine learning and multivariate statistical analysis based on the centering matrix, 2021.
- C. Heeyoul and C. Seungjin. Robust kernel Isomap, *ScienceDirect*, 2006, pp. 853-862.