

Chap 6 Classical Multidimensional Scaling.

New ~ Week 0

6.1 Intro to Multidim. Scaling.

Input: a similarity matrix or an array of similarity matrices.

Output: a low dimensional data set in some Euclidean space.

MDS = multidimensional scaling.

MDS $\left\{ \begin{array}{l} \text{qualitative (nonmetric MDS)} \\ \text{quantitative (metric MDS)} \end{array} \right.$

By number of similarity matrices,

- (1) **classical MDS**.
(one similarity matrix in an unweighted model).
- (2) **replicated MDS**.
(several .. matrices)
- (3) **weighted MDS**.
(.. .. . a weighted ..)
(i.e. different similarity matrices have different importance.)

6.2 Euclidean Metric and Gram Matrices.

In classical MDS (**CMDS**), the similarity matrix is assumed to come from data pts in an Euclidean space and distances b/w them.

Goal: Find "configuration" $Y \in \mathbb{R}^d$ s.t. the Euclidean distance matrix of Y best approximates the given similarities.

Def
 $X = \{x_1, \dots, x_n\} \in \mathbb{R}^D$. ($d_2 :=$ the Euclidean distance).
The **Euclidean distance matrix** on X is

$$D := [d_2(x_i, x_j)]_{i,j=1}^n$$

The **Euclidean square-distance matrix** on X is

$$S := [d_2^2(x_i, x_j)]_{i,j=1}^n$$

Rmk:

D and S above are symmetric and invariant of shift and rotation on X .

Notation:

For $a \in \mathbb{R}^D$, $X_a := \{x+a \mid x \in X\}$.

For rotation R (orthogonal w/ $\det 1$),

$$X_R := \{Rx \mid x \in X\}.$$

Def D : $n \times n$ symmetric matrix.

D : **Euclidean distance matrix** if

\exists integer $m > 0$, $\mathcal{I} = \{z_1, \dots, z_n\} \in \mathbb{R}^m$ s.t.

$$D = [d_2(z_i, z_j)]_{i,j=1}^n$$

Such \mathcal{I} is called a **configurative point set** (or a **configuration**) of D .

Q:

How to determine whether a given symmetric matrix is a Euclidean distance matrix?

Def $X = \{x_1, \dots, x_n\} \in \mathbb{R}^D$.

The **Gram matrix** on the data set X is

$$G_1 = [\langle x_i, x_j \rangle]_{i,j=1}^n$$

Rmk:

① $X := (x_1 \dots x_n)$, the $D \times n$ data matrix

Then $G_1 = X'X$. Thus, G_1 : **psd**. (positive semidef.)

② $D_{ij} = \sqrt{G_{ii} + G_{jj} - 2G_{ij}}$. **i.e.** $S = \delta(G_1)1^T + 1\delta(G_1)^T - 2G_1$.
where $\delta(G_1)$: diagonal elements of G_1 .
 $1 = [1, 1, \dots, 1]^T$.
i.e. we can recover S (so D) via G_1 .

Def $X = \{x_1, \dots, x_n\} \in \mathbb{R}^D$, a data set.

$$\bar{x} := \frac{1}{n} \sum_i x_i$$

$\hat{X} := X_{\bar{x}} = \{x_i - \bar{x} \mid i=1, \dots, n\}$, called a **centered data set**.

$\hat{X} := [\hat{x}_1 \dots \hat{x}_n]$, called the **centered data matrix**.

$$G_1^c = [\langle \hat{x}_i, \hat{x}_j \rangle]_{i,j=1}^n = \hat{X}'\hat{X}, \text{ called the}$$

In general, **centering Gram matrix** of X .

for $a \in \mathbb{R}^D$, the **a -shifted Gram matrix** of X

is $G_1^a = [\langle x_i - a, x_j - a \rangle]_{i,j=1}^n$.

Rmk:

As before, $D_{ij} = \sqrt{G_{ii}^c + G_{jj}^c - 2G_{ij}^c}$.

Def

$\mathbf{1} := [1, 1, \dots, 1]' \in \mathbb{R}^n$.

$\mathbf{E} := \mathbf{1}\mathbf{1}'$, the $n \times n$ matrix w/ all entries = 1.

$\mathbf{H} := \mathbf{I} - \frac{1}{n} \mathbf{E}$, called the n-centralizing matrix.

(If n is understood, omit "n-").

Lemma 6.1

\mathbf{H} has the following properties:

- (1) $\mathbf{H}^2 = \mathbf{H}$.
- (2) $\mathbf{1}'\mathbf{H} = \mathbf{H}\mathbf{1} = \mathbf{0}$.
- (3) \mathcal{X} : centered data set $\Leftrightarrow \mathbf{X}\mathbf{H} = \mathbf{X}$.
- (4) A psd \mathbf{C} is a centering Gram matrix $\Leftrightarrow \mathbf{H}\mathbf{C}\mathbf{H} = \mathbf{C}$.

<Pf>

- (1), (2): direct computation.
- (3): note \mathcal{X} : centered $\Leftrightarrow \mathbf{X}\mathbf{1} = \mathbf{0}$. (*)
- (4): (\Rightarrow) by (3). (\Leftarrow) use (*). (#)

Lemma 6.2

\mathbf{X} : data matrix, \mathbf{G} : Gram matrix of \mathbf{X} .

Then

$\mathbf{X}\mathbf{H}$: the centered data matrix of \mathbf{X} , and
 $\mathbf{H}\mathbf{G}\mathbf{H} = \mathbf{G}^c$, the centering Gram matrix of \mathbf{X} .

<Pf>

$\mathbf{X}\mathbf{H}$: by def. $\mathbf{H}\mathbf{G}\mathbf{H}$: by def. (#)

Def

\mathbf{A} : symmetric $n \times n$ matrix (not necessarily psd).

The centering matrix of \mathbf{A} , denoted \mathbf{A}^c , is

$\mathbf{A}^c = \mathbf{H}\mathbf{A}\mathbf{H}$. Rmk: A symmetric \mathbf{S} is centering $\Leftrightarrow \mathbf{H}\mathbf{S}\mathbf{H} = \mathbf{S}$.

Thm 6.1 (*) \mathcal{X} : data set.

$\mathbf{S} :=$ Euclidean square-distance matrix of \mathcal{X} .

$\mathbf{G}^c :=$ centering Gram matrix of \mathcal{X} .

Then

$\mathbf{G}^c = -\frac{1}{2} \mathbf{S}^c (= -\frac{1}{2} \mathbf{H}\mathbf{S}\mathbf{H})$.

(we can find \mathbf{G}^c via \mathbf{S}).

<Pf>

Straightforward computations.

(#)

The following gives an answer for how to determine whether a matrix is Euclidean distance:

Thm 6.2

\mathbf{A} : symmetric matrix. Then

- (1) \mathbf{A} : a Gram matrix $\Leftrightarrow \mathbf{A}$: psd.
- (2) \mathbf{A} : centering Gram $\Leftrightarrow \mathbf{A}$: centering psd.
- (3) \mathbf{A} : Euclidean square-distance matrix

$\Leftrightarrow -\frac{1}{2} \mathbf{A}^c$: centering psd, and $\mathbf{A} = \delta(\mathbf{G})\mathbf{1}\mathbf{1}' + \mathbf{1}\delta(\mathbf{G})' - 2\mathbf{G}$,
 where $\mathbf{G} = -\frac{1}{2} \mathbf{A}^c$.

6.3 Description of Classical Multidimensional Scaling.

Idea: (of metric multidimensional scaling (MDS)).

$\mathbf{D} = [d_{ij}]_{i,j=1}^n$ distance matrix of n object.

$d > 0$: an integer.

The goal is to find $\mathbf{Y} = \{y_1, \dots, y_n\} \in \mathbb{R}^d$ s.t.

$d_Y(y_i, y_j) \approx d_{ij}, \forall i, j. (*)$

In practice, closeness of (*) is measured by a reasonable "lost function".

Lemma 6.3

$\mathbf{D} = [d_{ij}]_{i,j=1}^n$, Euclidean distance matrix.

$\mathbf{S} = [d_{ij}^2]$: its square-distance

$\mathbf{G}^c := -\frac{1}{2} \mathbf{S}^c$ (i.e. its centering Gram matrix).

$r := \text{rank}(\mathbf{G}^c)$.

Then \exists r -dim. centered data set $\mathcal{X} = \{x_1, \dots, x_n\} \in \mathbb{R}^r$

s.t. $d_2(x_i, x_j) = d_{ij}, \forall i, j$.

Def

$r = \text{rank}(\mathbf{G}^c)$ in Lemma 6.3 is called the intrinsic configuration dimension of \mathbf{D} .

\mathcal{X} in L 6.3 is called an exact configurat'n of \mathbf{D} .

Rmk:

The lost function of CMDS is set to be

$\eta(\mathbf{Y}) = \sum_{i,j=1}^n (d_{ij}^2 - d_2^2(y_i, y_j))$, s.t. $\mathbf{Y} = \mathbf{T}(\mathcal{X})$, \mathbb{R}_{UI}^r

where \mathbf{T} : o.g. project'n from \mathbb{R}^r to some d -subsp. \mathbf{S}_d .

\mathcal{X} : an exact configurat'n of \mathbf{D} .

i.e.

we are solving the minimization problem:

$\arg \min \eta(\mathbf{Y})$, where $\mathbf{Y} = [y_1, \dots, y_n]$, $\mathbf{y} = \{y_1, \dots, y_n\}$, $\mathbf{Y} \in \mathcal{M}_{d,n}$

Lemma 6.4

$\mathcal{Z} = \{z_1, \dots, z_n\} \in \mathbb{R}^r$, a given data set.

$S_Z := [s_{ij}]$, where $s_{ij} = d_z^2(z_i, z_j)$.

G_Z^c := its centering Gram matrix.

Then $\text{tr}(G_Z^c) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^n s_{ij}$.

<Pf>

By Thm 6.1, $G_Z^c = \frac{1}{2} S_Z^c = \frac{1}{2} H S_Z H$.

$$\begin{aligned} \Rightarrow \text{tr}(G_Z^c) &= \text{tr}\left(\frac{1}{2} H S_Z H\right) = \text{tr}\left(\frac{1}{2} H^2 S_Z\right) = \frac{1}{2} \text{tr}(H S_Z) \\ &= \frac{1}{2} \text{tr}\left((I - \frac{1}{n} E) S_Z\right) = \frac{1}{2} [\text{tr}(S_Z) - \frac{1}{n} \text{tr}(E S_Z)] \\ &= \frac{1}{2n} \text{tr}(E S_Z) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^n s_{ij}. \quad (\#) \end{aligned}$$

Lemma 6.5

$D_Z := [d_z(z_i, z_j)]_{i,j=1}^n$, $\hat{Z} := [\hat{z}_1, \dots, \hat{z}_n]$.

Then $\|\hat{Z}\|_F = \frac{1}{\sqrt{2n}} \|D_Z\|_F$. (Recall: $\|A\|_F = \sqrt{\text{tr}(A^* A)}$)

<Pf>

$$\|D_Z\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n s_{ij} \stackrel{\text{by L6.4}}{=} 2n \text{tr}(G_Z^c) = 2n \cdot \text{tr}(\hat{Z}' \cdot \hat{Z})$$

$$= 2n \|\hat{Z}\|_F^2. \text{ Thus, } \|\hat{Z}\|_F = \frac{1}{\sqrt{2n}} \|D_Z\|_F. \quad (\#)$$

Thm 6.3

X : the centered configurat'n of D in Lemma 6.3.

SVD of X : $X = V \Sigma_r U'$, where (r : rank of X).

$$\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r), \quad r \times r, \quad \sigma_1 \geq \dots \geq \sigma_r > 0.$$

$$V = [v_1, \dots, v_r], \quad r \times r, \text{ orthogonal.}$$

$$U, \quad n \times r, \text{ w/ orthonormal columns.}$$

$$\text{For } d \leq r, \quad V_d := [v_1, \dots, v_d].$$

$$Y := V_d' X.$$

Then Y is a sol'n to the minimization problem stated before Lemma 6.4 w/

$$\eta(Y) = \sum_{i=1}^r \sigma_i^2.$$

Rmk:

Although the motivations of PCA and CMDS are different, when the input similarity matrix of CMDS is an Euclidean distance matrix, they are essentially

identical.

② CMDS can be recasted in terms of weight matrices: DR P3

The input similarity matrix is the weight matrix of n objects, say W .

The goal of CMDS is, under a given d , try to find $Y \in \mathbb{R}^d$, $Y = \{y_1, \dots, y_n\}$ s.t.

$$\|W - W(Y)\|_F^2 \text{ is minimized.}$$

this is exactly the loss function.

6.4 CMDS Algorithm.

Input: a similarity matrix D .

Step 1 Create centering Gram matrix.

$$G := \frac{1}{2} H (D \circ^2) H, \text{ where } D \circ^2: \text{pwise square of } D. \text{ (called the Hadamard square of } D).$$

Step 2 Make spectral decomp. of G .

$$r := \text{rank}(G).$$

$$G = U \Lambda U', \text{ spec. decomp. of } G, \text{ where}$$

$$U = [u_1, \dots, u_r], \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_r), \quad \lambda_1 \geq \dots \geq \lambda_r > 0.$$

Step 3 Find configurat'n.

$$U_d := [u_1, \dots, u_d]$$

$$\Sigma_d := \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_d})$$

$$\text{Then, the configurat'n is } Y = \Sigma_d U_d'.$$

Def

$$A = [a_{ij}], B = [b_{ij}], \text{ matrices of the same dim.}$$

The Hadamard product of A and B , denoted

$$A \boxtimes B, \text{ is } [a_{ij} b_{ij}]. \text{ (also called pwise product)}$$

Chap 7 Random Projection.

New Week

Problems of PCA: < 7.1 Inero.

PCA tries to minimize "global distortion" but may not preserve local separation of the data.

The following concept deals w/ local separati'n.

7.1.1 Lipschitz Embedding.

Def $f: \mathcal{X} \subseteq \mathbb{R}^D \rightarrow \mathbb{R}^k$ is called a Lipschitz embedding or Lipschitz mapping if \exists const. $A, B > 0$

$$\text{s.t. } A \|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq B \|u - v\|^2, \quad \forall u, v \in \mathcal{X}.$$

Rmk:

For general applicatn, we may use norms other than Euclidean norms.

7.1.2 J-L Embeddings.

Def $\mathcal{X} = \{x_1, \dots, x_n\} \in \mathbb{R}^D$.

$R = \{r_1, \dots, r_K\} \in \mathbb{R}^D$, K random vectors.

$X := [x_1 \dots x_n]$, $D \times n$.

$R := [r_1 \dots r_K]$, $D \times K$.

$Y := R'X = [y_1 \dots y_n]$, $K \times n$

The map $f: \mathcal{X} \rightarrow \mathcal{Y}$, $x_i \mapsto R'x_i$, is called a random projection.

Rmk:

① Random proj. proves a very useful tool in the construction of Lipschitz mappings.

② The feasibility of rand. proj. to construct Lip. map is due to a result by Johnson and Lindenstrauss. Thus, a Lip. map constructed as so is often called JL-embedding.

③ It is a kind of Monte-Carb method.

7.2 Random Proj. Algo.

7.2.1 Rand. Matrix and Rand. Proj.

Notation: r : r.v.

① $r \sim N(0,1)$ if r : standard normal.

② $r \sim U(\{-1,1\})$ if r : uniform on $\{-1,1\}$.

③ $r \sim U(\{\pm\sqrt{3}, 0\})$ if $r := \begin{cases} \sqrt{3} & \text{w/ prob. } 1/6 \\ 0 & \text{" " } 2/3 \\ -\sqrt{3} & \text{" " } 1/6 \end{cases}$

We shall call r.v. distributed as above of Type-1, Type-2, and Type-3, resp.

Rmk:

All of the three above have zero mean and unit variance.

Def

A random matrix is a matrix whose entries are iid r.v. Thus we can talk about random matrix of Type-1, Type-2, and Type-3.

Def

$R = [r_{ij}]$, a $K \times D$ random matrix of a type ^{certainly} P4.
A normalized random projection is a map $\mathbb{R}^D \rightarrow \mathbb{R}^K$ w/ $u \mapsto \frac{1}{\sqrt{K}} Ru$.

7.2.2 Rand. Proj. Algo.

2 steps:

(1) Random matrix creation.

(2) Matrix multiplication.

7.3 Justification.

7.3.1 Johnson and Lindenstrauss Lemma.

Lemma 7.1 (JL).

The statement is logically weird. !!!

Lemma 7.2

$R = [r_{ij}]$: $K \times D$ rand. matrix whose entries are iid w/ zero mean and unit var.

$f: \mathbb{R}^D \rightarrow \mathbb{R}^K$, $a \mapsto \frac{1}{\sqrt{K}} Ra$.

Then (1) $f(a)$: rand. vector in \mathbb{R}^K .

(a: fixed) (2) all components of $f(a)$ are iid w/ 0 mean and var. $\frac{1}{K} \|a\|^2$.

Thus, $E(\|f(a)\|^2) = \|a\|^2$.

Cor 7.1

f : as above.

$a \in \mathbb{R}^D$: unit.

Then $\sqrt{K} f(a)$: K -dim. rand. vector w/

iid entries w/ zero mean and unit var.
i.e. For $K \times D$ rand. matrix R w/ iid entries w/ zero mean and unit var. and unit vector a , Ra is a rand. vector whose entries are iid w/ zero mean and unit var.

7.3.2 Rand. Proj. based on Gaussian Dist.

Thm 7.1 (Dasgupta and Gupta) *

$0 < \epsilon < 1$, n : positive integer.

K : positive integer such that $K \geq 4(\epsilon^2/2 - \epsilon^2/3)^{-1} \ln n$.

Then, $\forall \mathcal{X} \in \mathbb{R}^D$, $|\mathcal{X}| = n$, \exists linear $f: \mathbb{R}^D \rightarrow \mathbb{R}^K$

s.t. $(1-\epsilon)\|u-v\|^2 \leq \|f(u)-f(v)\|^2 \leq (1+\epsilon)\|u-v\|^2$
 $\forall u, v \in \mathcal{X}$.

Def

An algorithm has randomized polynomial time if
 (1) it runs in poly. time in input size and
 (2) $NO \rightarrow NO$, and
 $YES \rightarrow YES$ w/ positive (const.) prob.

Lemma 7.3 ($K \in D$)

$R: K \times D$ rand. matrix of Gaussian type. (Type-1).
 $a \in \mathbb{R}^D$: unit vector.
 $y := Ra$, $\beta > 1$.

Then $\Pr[\|y\|^2 \leq K/\beta] < \exp(\frac{K}{2}(1-\frac{1}{\beta}-\ln\beta))$, and
 $\Pr[\|y\|^2 \geq K\beta] < \exp(\frac{K}{2}(1-\beta+\ln\beta))$.

Rank: (*) (Cor 7.2)

Using Lemma 7.3, we can prove:
 Under conditions of Thm 7.1, $f(a) := \frac{1}{\sqrt{K}} Ra$, then
 for $u \neq v$ in \mathcal{X} ,
 $\Pr(\frac{\|f(u)-f(v)\|^2}{\|u-v\|^2} \notin [1-\epsilon, 1+\epsilon]) < \frac{2}{n^2}$.

Denote the event $(\frac{\|f(u)-f(v)\|^2}{\|u-v\|^2} \notin [1-\epsilon, 1+\epsilon])$ by $A_{u,v}$.

Then
 $P((1-\epsilon)\|u-v\|^2 \leq \|f(u)-f(v)\|^2 \leq (1+\epsilon)\|u-v\|^2, \forall u,v)$
 $= P(\bigcap_{u \neq v} A_{u,v}^c) = 1 - P(\bigcup_{u \neq v} A_{u,v}) > 1 - \frac{n(n-1)}{2} \cdot \frac{2}{n^2}$
 $= 1/n > 0$.
 Thus, Thm 7.1 is proved and we also have:

Thm 7.1 (Continued.)

The mapping f in Thm 7.1 can be found in randomized polynomial time.

7.3.3 Rand. Proj. based on Type 2 and Type 3.
Thm 7.2 (Achlioptas).

Conditions as in Thm 7.1.
 $\mathcal{X} \subseteq \mathbb{R}^D$, $|\mathcal{X}| = n$, $X = [x_1 \dots x_n]$
 $R: K \times D$ rand matrix of Type 2 or Type 3.
 $f: \mathbb{R}^D \rightarrow \mathbb{R}^K$, $u \mapsto \frac{1}{\sqrt{K}} Ru$.
 Then f satisfies the following Lipschitz condition:
 $(1-\epsilon)\|u-v\|^2 \leq \|f(u)-f(v)\|^2 \leq (1+\epsilon)\|u-v\|^2, \forall u,v \in \mathcal{X}$,
 w/ prob. at least $1 - 1/n^\beta$, where $\beta = \frac{K(\epsilon^{2/4} - \epsilon^{3/6})}{\ln n} - 2 > 0$.

<PF>

Use conditions to conclude:

$$\Pr(\frac{\|f(u)-f(v)\|^2}{\|u-v\|^2} \notin [1-\epsilon, 1+\epsilon]) < \frac{2}{n^{2+\beta}} \quad (*)$$

Part III Nonlinear Dimensionality Reduction.

Chap 8 Isomaps.

8.1 Isomap Embeddings.

8.1.1 Description of Isomaps.

Motivation: use geodesic metric instead of Euclidean metric.

Method:

Measure dissimilarity using geodesic metric, which is approximated by graph distance.

Def

A map is called an isometric map if it preserves distances.

Rank:

① "Isomap" is an abbreviation of isometric mapping.

② This method was introduced by deSilva, et al.

Mathematical description:

Assume $\mathcal{X} = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^D$ lies on a d -dim. Riemannian manifold $M \subseteq \mathbb{R}^D$, w/ $d \ll D$.

$d_M :=$ geodesic metric on M .

Then \exists chart $f: M \rightarrow \mathbb{R}^d$ st.

$$d_2(f(x), f(z)) = d_M(x, z), \forall x, z \in M.$$

In particular, geodesic distances on \mathcal{X} is preserved by f .

8.1.2 Geodesic metric on discrete set.

Motivation:

Though geodesic metric is infeasible due to unknown formula of M , we can use graph distance to app. it.

Def

$\mathcal{X} = \{x_1, \dots, x_n\}$: a data set in \mathbb{R}^D .
 $G = [\mathcal{X}, E]$, a graph indicating the nbd system on \mathcal{X}
 Define $d_G: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, called the graph distance,
 by ① if $(x, y) \in E$, then $d_G(x, y) := d_2(x, y)$.

(iii) If $(x, y) \notin E$, for a path $\gamma = (x_0, x_1, \dots, x_{st+1})$

w/ $x_0 = x, x_{st+1} = y$, the path distance

$$d_\gamma(x, y) := d_2(x_0, x_1) + \dots + d_2(x_s, x_{st+1}).$$

$$d_G(x, y) := \min \{d_\gamma(x, y) \mid \gamma: \text{path from } x \text{ to } y\}.$$

Rank:

When the data set is dense enough on M , d_G app. d_M under certain conditions.

8.1.3 Isomap kernel and its Const. Shift.

$$D_G := [d_G(x_i, x_j)], \text{ where } \mathcal{X} = \{x_1, \dots, x_n\}.$$

Note that d_G is a true metric on \mathcal{X} .

case 1 \mathcal{X} is dense enough on M so that D_G app.

$$D_M = [d_M(x_i, x_j)] \text{ very well.}$$

Then, for a preassigned target dimension K , we may apply CMDS on D_G to find $\mathcal{Y} = \{y_1, \dots, y_n\} \subseteq \mathbb{R}^K$ w/ $[d_2(y_i, y_j)] \approx D_G \approx D_M$.

More precisely,

$$S := D_G \otimes D_G = [d_G^2(x_i, x_j)].$$

By Thm 6.1, $G^c = -\frac{1}{2} S^c = -\frac{1}{2} H S H$ is the centered Gram matrix.

Applying PCA, we can obtain \mathcal{Y} (depending on the preassigned K). G^c is called the Isomap kernel.

case 2 \mathcal{X} is not dense enough on M .

Issue: D_G not app. D_M well.

$\Rightarrow G^c$ may not be positive semi-definite.

\Rightarrow CMDS not applicable.

Simple Sol'n: add a $\tau > 0$ as a compensation.

$$d_\tau(x_i, x_j) := \begin{cases} d_G(x_i, x_j) + \tau, & i \neq j. \\ 0, & i = j. \end{cases} \quad \begin{matrix} \text{(constant-shift)} \\ \text{technique} \end{matrix}$$

Thus, $S_\tau := [d_\tau^2(x_i, x_j)]$ and the centered Gram

$$\text{is } G_\tau^c := -\frac{1}{2} S_\tau^c.$$

G_τ^c is called a constant-shifted Isomap kernel.

claim: $\exists \tau > 0$ s.t. G_τ^c psd.

Thm 8.1 (Cailliez).

The minimal $\tau = \tau^*$ making G_τ^c psd is the largest eigenvalue of

$$B = \begin{pmatrix} 0 & -S^c \\ -I & 2D^c \end{pmatrix}.$$

8.2 Isomap Algorithm.

8.2.1 Algorithm Description

Step 1 Neighborhood Definition.

$\mathcal{X} = \{x_1, \dots, x_n\}$: a pt set in \mathbb{R}^d .

Use either K -neighborhood or ϵ -neighborhood to form a neighborhood system on \mathcal{X} .

The method is called K -Isomap or ϵ -Isomap resp.

Step 2 Graph distance computation.

Denote the graph created via the nbd system by $G = [\mathcal{X}, E]$.

case 1 ϵ -nbd.

Compute the graph distance matrix.

case 2 K -nbd.

Since the graph is not simple, to ensure the matrix to be symmetric, (1) either define

$$d_G(x_i, x_j) = \min \{d_G(i, j), d_G(j, i)\}, \text{ or}$$

(2) modify G to be undirected.

Step 3 DR kernel construction.

$$S_G := [d_G^2(i, j)].$$

$$G^c := -\frac{1}{2} H S_G H = -\frac{1}{2} S_G^c.$$

case 1 G^c : psd, done.

case 2 G^c : not psd. Use $G_{\tau^*}^c$ instead, where

$$\tau^* = \text{maximum eigenvalue of } \begin{pmatrix} 0 & -S_G^c \\ -I & 2D_G^c \end{pmatrix}.$$

Step 4 CMDS.

8.3 Dijkstra's algorithm.

For dense graph, use Floyd's algorithm.

For sparse " " " " Dijkstra's " "

See P160-P169.

Chap 9 Maximum Variance Unfolding.

Motivation:

Preserving both local distances and angles.

AKA:

(1) MVU (maximum variance unfolding)

(2) SDE (semidefinite embedding)

(reason: involving semidefinite programming (SDP)).

Efficiency Modification:

landmark MVU (LMVU).

9.1 MVU Method Description.

9.1.1 Description of MVU Method.

$\mathcal{X} \subseteq \mathbb{R}^D$, lies on a convex d-mfd M .

h : an isometric chart on (the whole) M .

$\mathcal{Y} := h(\mathcal{X}) \subseteq \mathbb{R}^d$, which is the target.

Suppose we have built a nbd system on \mathcal{X} .

For $x_i \in \mathcal{X}$, let P_i be the nbd of x_i .

Let $Q_i := h(P_i)$. $y_i := h(x_i)$. $N(i) := \{j \mid x_j \in P_i\}$.

Idea: $\|x_j - x_k\| \approx \|y_j - y_k\|$, $\forall j, k \in N(i)$.

Thus we want to obtain:

$$y = \operatorname{argmin}_{y \in \mathbb{R}^d} \sum_{i=1}^n \sum_{j \sim k} (\|x_k - x_j\|^2 - \|y_k - y_j\|^2)^2, \text{ where}$$

$j \sim k, \forall j, k \in N(i), \text{ including } i$.

Issue:

The above optimization problem is nonconvex; hard to solve directly.

Sol'n:

Consider "another similar optimization problem" solvable by SDP.

Let $Y := (y_1, \dots, y_n)$, where $\mathcal{Y} = \{y_1, \dots, y_n\}$.

$K := Y^T Y$, the Gram matrix.

WLOG, we may assume \mathcal{Y} is centered.

i.e. $y_1 + \dots + y_n = 0$. i.e. $\sum_{i,j=1}^n K_{ij} = 0$. (X1) centering constraint.

$$S_{kj} := \|y_j - y_k\|^2.$$

Then

$$S_{kj} = K_{kk} + K_{jj} - 2K_{kj}, \forall j, k \in N(i), 1 \leq i \leq n. \quad (X2)$$

(We may assume K_{ij} to eliminate redundancy) local geometry constraint

Note that (X2) is a constraint on local geometry.
(X1) " " convenient constraint.

MVU, as the name suggests, has the following constraint:

$$y = \operatorname{argmax}_{y \in \mathbb{R}^d} \sum_{i,j=1}^n \|y_i - y_j\|^2 \quad (X3) \leftarrow \text{pull } y \text{ as far apart as possible.}$$

$w/(X1), (X2)$

We want everything to be stated in terms of K , so we make the following deduction:

By Thm 6.1, $G_Y^c = -\frac{1}{2} S_Y^c = -\frac{1}{2} H S_Y H$.

Note that $K = G_Y^c$ (due to (X1)).

$$\Rightarrow \operatorname{tr}(K) = \operatorname{tr}(G_Y^c) = -\frac{1}{2} \operatorname{tr}(H S_Y H)$$

$$\stackrel{(H^2=I)}{\cong} -\frac{1}{2} \operatorname{tr}(S_Y H) = -\frac{1}{2} \operatorname{tr}(S_Y (I - \frac{1}{n} E))$$

$$\stackrel{(\operatorname{tr}(S_Y)=0)}{\cong} \frac{1}{2n} \operatorname{tr}(S_Y E) = \frac{1}{2n} \sum_{i,j=1}^n \|y_i - y_j\|^2.$$

Thus, (X3) can be rephrased as

$$y = \operatorname{argmax}_{y \in \mathbb{R}^d} \operatorname{tr}(K).$$

$w/(X1), (X2)$

In summary, we obtain the kernel K via

$$\begin{aligned} &\text{maximize } \operatorname{tr}(K), \quad K: n \times n \text{ psd (i.e. } K \geq 0). \\ &\text{s.t. } \sum_{i,j=1}^n K_{ij} = 0, \\ &\quad K_{kk} + K_{jj} - 2K_{kj} = S_{kj}, \quad k, j \in N(i) \\ &\quad \quad \quad 1 \leq i \leq n. \end{aligned}$$

Remark:

① The input can either be (i) \mathcal{X} or (ii) D_X (equivalently, S_X).

② (X2) requires the neighbor relations to be complete. i.e. for each i , x_j and x_k are nb, $\forall j, k \in N(i)$. We can achieve this by locally completing the graph.

9.1.2 MVU Algorithm.

Assume $X \in \mathbb{R}^D$ or a local dist. matrix D is given.

Step 1. Nbd definition. (~~#~~ ^{when} only X is given).

k -nbd is commonly used. $A :=$ adjacency matrix.

Step 2. Compute S or G^c .

Step 3. Constraints Generating.

If $A(k, j) = 1$, set the constraint

$$K_{kk} + K_{jj} - 2K_{kj} = s_{kj} \text{ (or } G_{kk} + G_{jj} - 2G_{kj}).$$

Also, we need to add in the centering constraint

$$\sum_{k,j=1}^n K_{kj} = 0.$$

Step 4. Build MVU kernel using SDP.

Apply SDP to obtain K .

Step 5. Spectral Decomp. of K .

Choose a d and find out the d -leading eigenvec. of K .

9.2 SDP.

Def

For $n \times n$ sym. matrices A and B , define their

inner product $A \cdot B$ by $A \cdot B = \sum_{i=1}^n \sum_{j=1}^n A_{ij} B_{ij}$.

Given $C, A_i, i=1, \dots, m$, symmetric $n \times n$.

$b_i \in \mathbb{R}, i=1, \dots, m$.

The following optimization problem is called a

semidefinite programming. (SDP).

minimize $C \cdot X$

s.t. $A_i \cdot X = b_i, i=1, \dots, m$.

$X \succeq 0$.

Def

The set of $n \times n$ sym. matrices satisfying the constraints is called the **feasible set**.

Prop

An SDP has a soln if the feasible set $\neq \emptyset$.

Rmk:

MVU can be solved by SDP.