Big Data Analysis (MA60306)

Bibhas Adhikari

Spring 2022-23, IIT Kharagpur

Lecture 12 February 3, 2023

3

A D N A B N A B N A B N

 $\rightarrow\,$ Thus in an MSD processing the input is a similarity matrix and the out put is a low dimensional space, which is usually 2 or 3 dimensional

★ ∃ ► < ∃ ►</p>

 $\rightarrow\,$ Thus in an MSD processing the input is a similarity matrix and the out put is a low dimensional space, which is usually 2 or 3 dimensional

Classical MDS In this case, we assume that the data set is

$$\mathcal{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\}, \mathbf{x}_j \in \mathbb{R}^d,$$

and the data matrix is

$$X = [\mathbf{x}_1, \ldots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$$

Note that each column represents a data point here.

 $\rightarrow\,$ Thus in an MSD processing the input is a similarity matrix and the out put is a low dimensional space, which is usually 2 or 3 dimensional

Classical MDS In this case, we assume that the data set is

$$\mathcal{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\}, \mathbf{x}_j \in \mathbb{R}^d,$$

and the data matrix is

$$X = [\mathbf{x}_1, \ldots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$$

Note that each column represents a data point here. The Euclidean distance matrix is $D = [d_{ii}] = [d_2(\mathbf{x}_i, \mathbf{x}_i)]$ where

$$d_2(\boldsymbol{x}, \boldsymbol{y}) = \sqrt{\sum_{j=1}^n (x_j - y_j)^2}$$

We define Euclidean square-distance matrix

$$S = [d_2^2(\boldsymbol{x}_i, \boldsymbol{x}_j)]$$

2

A D N A B N A B N A B N

We define Euclidean square-distance matrix

$$S = [d_2^2(\boldsymbol{x}_i, \boldsymbol{x}_j)]$$

Then observe that D and S are \triangle Symmetric

3

A D N A B N A B N A B N

We define Euclidean square-distance matrix

$$S = [d_2^2(\boldsymbol{x}_i, \boldsymbol{x}_j)]$$

Then observe that D and S are

 \triangle Symmetric

 \bigtriangleup invariant under shift and rotation

э

・ 何 ト ・ ヨ ト ・ ヨ ト

We define Euclidean square-distance matrix

$$S = [d_2^2(\boldsymbol{x}_i, \boldsymbol{x}_j)]$$

Then observe that D and S are

 \triangle Symmetric

 \bigtriangleup invariant under shift and rotation

Euclidean distance metric Any symmetric matrix D is called a Euclidean distance matrix or Euclidean metric if there exists a positive integer k and a set $\mathcal{Z} = \{z_1, \ldots, z_n\} \subset \mathbb{R}^k$ such that

$$D = [d_2(\boldsymbol{z}_i, \boldsymbol{z}_j)].$$

In that case \mathcal{Z} is called a *configuration* of D.

Question Can a matrix be detected as a Euclidean metric from its properties without finding the configuration explicitly?

3

< □ > < □ > < □ > < □ > < □ > < □ >

Question Can a matrix be detected as a Euclidean metric from its properties without finding the configuration explicitly?

Gram (Gramian) matrix of a data set Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$. Then the Gram matrix of \mathcal{X} is defined by

$$G = [g_{ij}] = [\mathbf{x}_i^T \mathbf{x}_j] = [\langle \mathbf{x}_i, \mathbf{x}_j \rangle]$$

3

< □ > < □ > < □ > < □ > < □ > < □ >

Question Can a matrix be detected as a Euclidean metric from its properties without finding the configuration explicitly?

Gram (Gramian) matrix of a data set Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$. Then the Gram matrix of \mathcal{X} is defined by

$$G = [g_{ij}] = [\mathbf{x}_i^T \mathbf{x}_j] = [\langle \mathbf{x}_i, \mathbf{x}_j \rangle]$$

Then G is a positive semi-definite matrix.

・ 何 ト ・ ヨ ト ・ ヨ ト

Question Can a matrix be detected as a Euclidean metric from its properties without finding the configuration explicitly?

Gram (Gramian) matrix of a data set Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$. Then the Gram matrix of \mathcal{X} is defined by

$$G = [g_{ij}] = [\mathbf{x}_i^T \mathbf{x}_j] = [\langle \mathbf{x}_i, \mathbf{x}_j \rangle]$$

Then G is a positive semi-definite matrix. However, every positive semi-definite matrix G has Cholesky decomposition i.e.

$$G = R^T R$$

for some $R = [\mathbf{r}_1, \ldots \mathbf{r}_n]$

<日

<</p>

Question Can a matrix be detected as a Euclidean metric from its properties without finding the configuration explicitly?

Gram (Gramian) matrix of a data set Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$. Then the Gram matrix of \mathcal{X} is defined by

$$G = [g_{ij}] = [\mathbf{x}_i^T \mathbf{x}_j] = [\langle \mathbf{x}_i, \mathbf{x}_j \rangle]$$

Then G is a positive semi-definite matrix. However, every positive semi-definite matrix G has Cholesky decomposition i.e.

$$G = R^T R$$

for some $R = [\mathbf{r}_1, \dots \mathbf{r}_n]$ Question What is the conclusion?

<日

<</p>

Note that for any $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$,

$$d_2(\mathbf{x}, \mathbf{y}) = \sqrt{\langle \mathbf{x}, \mathbf{x}
angle + \langle \mathbf{y}, \mathbf{y}
angle - 2 \langle \mathbf{x}, \mathbf{y}
angle}$$

i.e.

$$d_{ij}=d_2(\pmb{x}_i,\pmb{x}_j)=\sqrt{g_{ii}+g_{jj}-2g_{ij}}$$

3

<ロト < 四ト < 三ト < 三ト

Note that for any $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d$,

$$d_2(\mathbf{x},\mathbf{y}) = \sqrt{\langle \mathbf{x},\mathbf{x}
angle + \langle \mathbf{y},\mathbf{y}
angle - 2 \langle \mathbf{x},\mathbf{y}
angle}$$

i.e.

$$d_{ij}=d_2(oldsymbol{x}_i,oldsymbol{x}_j)=\sqrt{g_{ii}+g_{jj}-2g_{ij}}$$

If the given set of data points $\mathcal{X} = \{x_1, \dots, x_n\}, x_j \in \mathbb{R}^d$ lie on a k dimensional affine subspace (hyperplane) $H \subset \mathbb{R}^d$ then the center of

$$\overline{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$$

lies in *H*, and $S = H - \overline{x} = \{x - \overline{x} : x \in H\} \subset \mathbb{R}^d$ is a *d*-dimensional subspace parallel to *H*

Set $\widehat{\mathcal{X}} = \{\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_n\}, \ \widehat{\mathbf{x}}_i = \mathbf{x}_i - \overline{\mathbf{x}}$ is called the *centered data set* and the corresponding data matrix \widehat{X} is called the centered data matrix.

(B)

Set $\widehat{\mathcal{X}} = \{\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_n\}, \ \widehat{\mathbf{x}}_i = \mathbf{x}_i - \overline{\mathbf{x}}$ is called the *centered data set* and the corresponding data matrix $\widehat{\mathcal{X}}$ is called the centered data matrix. The Gram matrix corresponding to $\widehat{\mathcal{X}}$ given by

$$G^{c} = [\langle \widehat{\boldsymbol{x}}^{T}, \widehat{\boldsymbol{x}}_{j} \rangle] = \widehat{X}^{T} \widehat{X}$$

is called the *centering Gram matrix* of \mathcal{X} .

Set $\widehat{\mathcal{X}} = \{\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_n\}, \ \widehat{\mathbf{x}}_i = \mathbf{x}_i - \overline{\mathbf{x}}$ is called the *centered data set* and the corresponding data matrix $\widehat{\mathcal{X}}$ is called the centered data matrix. The Gram matrix corresponding to $\widehat{\mathcal{X}}$ given by

$$G^{c} = [\langle \widehat{\boldsymbol{x}}^{T}, \widehat{\boldsymbol{x}}_{j} \rangle] = \widehat{X}^{T} \widehat{X}$$

is called the *centering Gram matrix* of \mathcal{X} . If $G^{c} = [g_{ij}^{c}]$ then

$$d_{ij}=\sqrt{g^c_{ij}+g^c_{jj}-2g^c_{ij}}$$

Set $\widehat{\mathcal{X}} = \{\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_n\}, \ \widehat{\mathbf{x}}_i = \mathbf{x}_i - \overline{\mathbf{x}}$ is called the *centered data set* and the corresponding data matrix $\widehat{\mathcal{X}}$ is called the centered data matrix. The Gram matrix corresponding to $\widehat{\mathcal{X}}$ given by

$$G^{c} = [\langle \widehat{\boldsymbol{x}}^{T}, \widehat{\boldsymbol{x}}_{j} \rangle] = \widehat{X}^{T} \widehat{X}$$

is called the *centering Gram matrix* of \mathcal{X} . If $G^{c} = [g_{ij}^{c}]$ then

$$d_{ij}=\sqrt{g^{\,c}_{ij}+g^{\,c}_{jj}-2g^{\,c}_{ij}}$$

(revisiting) Centering/centralizing matrix Let $\mathbf{1} = [1, 1, ..., 1]^T \in \mathbb{R}^n$. Let $E = \mathbf{1}\mathbf{1}^T$. Then $C_n = I_n - \frac{1}{n}E$ Then

Then

$$\rightarrow C_n^2 = C_n$$

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣

Then

$$\begin{array}{l} \rightarrow \ \ C_n^2 = C_n \\ \rightarrow \ \ \mathbf{1}^T C_n = C_n \mathbf{1} = 0 \end{array}$$

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ― 臣

Then

$$\begin{array}{l} \rightarrow \ C_n^2 = C_n \\ \rightarrow \ \mathbf{1}^T C_n = C_n \mathbf{1} = 0 \end{array}$$

→ A data set $X = \{x_1, ..., x_n\}$ is centered if and only if $XC_n = X$, where $X = [x_1, ..., x_n]$

3

< □ > < 同 > < 回 > < 回 > < 回 >

Then

- $\rightarrow C_n^2 = C_n$
- $\rightarrow \mathbf{1}^T C_n = C_n \mathbf{1} = \mathbf{0}$
- $\label{eq:constraint} \begin{array}{l} \rightarrow \mbox{ A data set } \mathcal{X} = \{ \pmb{x}_1, \dots, \pmb{x}_n \} \mbox{ is centered if and only if } XC_n = X, \\ \mbox{ where } X = [\pmb{x}_1, \dots, \pmb{x}_n] \end{array}$
- \rightarrow A positive semi-definite matrix M is a centering Gram matrix if and only if $C_nMC_n = M$

э

A B A A B A

Then

- $\rightarrow C_n^2 = C_n$
- $\rightarrow \mathbf{1}^T C_n = C_n \mathbf{1} = \mathbf{0}$
- $\rightarrow A \text{ data set } \mathcal{X} = \{ \mathbf{x}_1, \dots, \mathbf{x}_n \} \text{ is centered if and only if } XC_n = X, \\ \text{where } X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$
- \rightarrow A positive semi-definite matrix M is a centering Gram matrix if and only if $C_nMC_n = M$

Conclusion Let G denote the Gram matrix of a data matrix X. Then the data matrix corresponding to the centered data set of X is XC_n , and the centering Gram matrix of X is $G^c = C_n GC_n$

Note that the data points are columns of X

3

A B A A B A

Theorem Let \mathcal{X} be a data set. Then

$$G^c = -\frac{1}{2}S^c$$

Proof Note that $\sum_{i=1}^{n} g_{ij}^{c} = 0$.

3

A D N A B N A B N A B N

Theorem Let \mathcal{X} be a data set. Then

$$G^c = -\frac{1}{2}S^c$$

Proof Note that $\sum_{i=1}^{n} g_{ij}^{c} = 0$. Then $d_{ij} = \sqrt{g_{ij}^{c} + g_{jj}^{c} - 2g_{ij}^{c}}$ implies

$$\sum_{i=1}^{n} d_{ij}^{2} = ng_{jj}^{c} + \sum_{i=1}^{n} g_{ii}^{c} \text{ and } \sum_{j=1}^{n} d_{ij}^{2} = ng_{ii}^{c} + \sum_{j=1}^{n} g_{jj}^{c}.$$

- 3

イロト 不得下 イヨト イヨト

Theorem Let \mathcal{X} be a data set. Then

$$G^c = -\frac{1}{2}S^c$$

Proof Note that $\sum_{i=1}^{n} g_{ij}^{c} = 0$. Then $d_{ij} = \sqrt{g_{ij}^{c} + g_{jj}^{c} - 2g_{ij}^{c}}$ implies

$$\sum_{i=1}^{n} d_{ij}^{2} = ng_{jj}^{c} + \sum_{i=1}^{n} g_{ii}^{c} \text{ and } \sum_{j=1}^{n} d_{ij}^{2} = ng_{ii}^{c} + \sum_{j=1}^{n} g_{jj}^{c}.$$

Thus

$$\begin{split} [S^c]_{ij} &= D_{ij}^2 - \frac{1}{n} \left(\sum_{i=1}^n d_{ij}^2 + \sum_{j=1}^n d_{ij}^2 - \frac{1}{n} \sum_{i,j=1}^n d_{ij}^2 \right) \\ &= d_{ij}^2 - g_{ii}^c - g_{jj}^c = -2g_{ij}^c \end{split}$$

3

イロト 不得下 イヨト イヨト

Theorem Let \mathcal{X} be a data set. Then

$$G^c = -\frac{1}{2}S^c$$

Proof Note that $\sum_{i=1}^{n} g_{ij}^{c} = 0$. Then $d_{ij} = \sqrt{g_{ij}^{c} + g_{jj}^{c} - 2g_{ij}^{c}}$ implies

$$\sum_{i=1}^{n} d_{ij}^{2} = ng_{jj}^{c} + \sum_{i=1}^{n} g_{ii}^{c} \text{ and } \sum_{j=1}^{n} d_{ij}^{2} = ng_{ii}^{c} + \sum_{j=1}^{n} g_{jj}^{c}.$$

Thus

$$\begin{split} [S^c]_{ij} &= D_{ij}^2 - \frac{1}{n} \left(\sum_{i=1}^n d_{ij}^2 + \sum_{j=1}^n d_{ij}^2 - \frac{1}{n} \sum_{i,j=1}^n d_{ij}^2 \right) \\ &= d_{ij}^2 - g_{ii}^c - g_{jj}^c = -2g_{ij}^c \end{split}$$

Question What does this equality mean when the data points are normalized?

Bibhas Adhikari (Spring 2022-23, IIT Kharag_l

The we have the following consequence: A matrix A is a Euclidean square-distance matrix if and only if $-\frac{1}{2}A^c$ is a centering positive semi-definite matrix.

3

・ 何 ト ・ ヨ ト ・ ヨ ト

The we have the following consequence: A matrix A is a Euclidean square-distance matrix if and only if $-\frac{1}{2}A^c$ is a centering positive semi-definite matrix.

Classical multidimensional scaling method Let $D = [d_{ij}]$ be a given distance matrix for a set of *n* objects. The we want to find a configuration $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$ such that a certain distance matrix associated with \mathcal{X} is as close as possible to D, i.e.

 $d_X(\boldsymbol{x}_i, \boldsymbol{x}_j) \approx d_{ij}$

for $1 \leq i, j \leq n$.

Lemma Suppose $D = [d_{ij}]$ is an $n \times n$ Euclidean metric and $S = [d_{ij}^2]$ is the corresponding square-distance matrix. Let $G^c = -\frac{1}{2}S^c$. If the rank of G^c is k then there is an k-dimensional centered vector set $\mathcal{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\} \subset \mathbb{R}^k$ such that

$$d_2(\boldsymbol{x}_i, \boldsymbol{x}_j) = d_{ij}, \ 1 \leq i, j \leq n.$$

< 回 > < 回 > < 回 >

Lemma Suppose $D = [d_{ij}]$ is an $n \times n$ Euclidean metric and $S = [d_{ij}^2]$ is the corresponding square-distance matrix. Let $G^c = -\frac{1}{2}S^c$. If the rank of G^c is k then there is an k-dimensional centered vector set $\mathcal{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\} \subset \mathbb{R}^k$ such that

$$d_2(\boldsymbol{x}_i, \boldsymbol{x}_j) = d_{ij}, \ 1 \leq i, j \leq n.$$

Proof By above G^c is a centering gram matrix. If rank of G^c is k then $G^c = X^T X$ for some matrix $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{k \times n}$, which has the desired property.

Lemma Suppose $D = [d_{ij}]$ is an $n \times n$ Euclidean metric and $S = [d_{ij}^2]$ is the corresponding square-distance matrix. Let $G^c = -\frac{1}{2}S^c$. If the rank of G^c is k then there is an k-dimensional centered vector set $\mathcal{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\} \subset \mathbb{R}^k$ such that

$$d_2(\boldsymbol{x}_i, \boldsymbol{x}_j) = d_{ij}, \ 1 \leq i, j \leq n.$$

Proof By above G^c is a centering gram matrix. If rank of G^c is k then $G^c = X^T X$ for some matrix $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{k \times n}$, which has the desired property.

Let us call k as the configuration of D and \mathcal{X} as the exact configuration of D. However, note that k can be large to meet the goal that we have low-dimensional configuration, say $d \ll k$.

- 小田 ト イヨト 一日

Question Can we take help of the PCA?

2

A D N A B N A B N A B N

Question Can we take help of the PCA?

Let Y denote a random desired matrix. Then consider the loss function

$$\mathcal{L}(\mathcal{Y}) = \sum_{i,j=1}^{n} \left(d_{ij}^2 - d_2^2(\boldsymbol{y}_i, \boldsymbol{y}_j) \right)$$

where \mathcal{Y} is obtained by orthogonal projection from \mathbb{R}^k to a *d*-dimensional subspace of \mathbb{R}^k and \mathcal{X} is the exact configuration of *D*.

Lemma Let $\mathcal{Z} \subset \mathbb{R}^k$ be a given data with Euclidean square-distance matrix $S_Z = [s_{ij}]$, where $s_{ij} = d_2^2(\mathbf{z}_i, \mathbf{z}_j)$, and let G_Z^c be its centering Gram matrix. Then

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

Proof Homework

3

・ 何 ト ・ ヨ ト ・ ヨ ト

Lemma Let $Z \subset \mathbb{R}^k$ be a given data with Euclidean square-distance matrix $S_Z = [s_{ij}]$, where $s_{ij} = d_2^2(\mathbf{z}_i, \mathbf{z}_j)$, and let G_Z^c be its centering Gram matrix. Then

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

Proof Homework

Lemma Let $D_Z = [d_2(z_i, z_j)]$ and $\widehat{Z} = [\widehat{z}_1, \dots, \widehat{z}_n]$, the centered data matrix corresponding to \mathcal{Z} . Then

$$\|\widehat{Z}\|_F = \frac{1}{\sqrt{2n}} \|D_Z\|_F$$

Proof Homework

医静脉 医原体 医原体 医原

Theorem¹ Let $\mathcal{X} \subset \mathbb{R}^k$ be the configuration of D such that \mathcal{X} is centered and the SVD of X be given by

$$X = U \Sigma_k V^T,$$

where $\Sigma_k = \text{diag}(\sigma_1, \ldots, \sigma_k)$, $U = [\boldsymbol{u}_1, \boldsymbol{u}_2, \ldots, \boldsymbol{u}_k]$. For a given $d \ll k$, let $U_d = [\boldsymbol{u}_1, \boldsymbol{u}_2, \ldots, \boldsymbol{u}_d]$ and $Y = U_d^T X$. Then Y is a solution of the above optimization problem with

$$\mathcal{L}(\mathcal{Y}) = \sum_{i=d+1}^{k} \sigma_i^2.$$

Proof Homework

¹Chapter 6, J. Wang, Geometric Structure of High-Dimensional Data and Dimensionality Reduction, Springer, 2012

Bibhas Adhikari (Spring 2022-23, IIT Kharag

Big Data Analysis

The classical MDS algoritm:

Step 1 Let *D* be the given distance matrix. Set $G^c = -\frac{1}{2}S^c$

- 20

イロト 不得下 イヨト イヨト

The classical MDS algoritm:

Step 1 Let *D* be the given distance matrix. Set $G^c = -\frac{1}{2}S^c$ **Step 2** Suppose rank of G^c is *k*. Compute the spectral decomposition of G^c as $G^c = U\Lambda U^T$, where $U = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_k]$, and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_k)$ with $\lambda_1 \ge \lambda_2 \le \dots \le \lambda_k$. The classical MDS algoritm:

Step 1 Let *D* be the given distance matrix. Set $G^c = -\frac{1}{2}S^c$ **Step 2** Suppose rank of G^c is *k*. Compute the spectral decomposition of G^c as $G^c = U\Lambda U^T$, where $U = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_k]$, and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_k)$ with $\lambda_1 \ge \lambda_2 \le \dots \le \lambda_k$. **Step 3** Set $U_d = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_k]$ and $\Sigma_d = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_d})$. Then the configuration is $Y = \Sigma_d U_d^T$

```
import numpy as np
from numpy.linalg import eig
D = np.array([[0,4,3,7,8],[4,0,1,6,7],[3,1,0,5,7],[7,6,5,0,1], [8,7,7,1,0]])
D2 = np.square(D)
C = np.eye(5) - 0.2*np.ones(5)
M = -0.5* C @ D2 @ C
I,V = eig(M)
s = np.real(np.power(1,0.5))
V2 = V[:,[0,1]]
s2 = np.diag(s[0:2])
Q = V2 @ s2
import matplotlib.pyplot as plt
plt.plot(Q[:,0],Q[:,1],'ro')
plt.show()
```

- 31

A (10) < A (10) < A (10) </p>