

* Fundamentals of BigData Analytics / 24. Nov. 2017

4.2 Multi-dimensional Scaling

$$p_1, \dots, p_n$$

$$\delta_{ij} = \delta_{ji} \neq 0, \quad \forall (i,j) \quad ; \quad \delta_{ii} = 0$$

$$\Delta = (\delta_{ij})_{i,j} \quad \text{dissimilarity matrix}$$

$$\mathcal{U}_n = \left\{ \Delta = (\delta_{ij})_{\substack{i,j \in n \\ i < j}} \mid \delta_{ij} = \delta_{ji} \geq 0, \delta_{ii} = 0 \quad \forall (i,j) = 1, \dots, n \right\}$$

$$x_1, \dots, x_n \in \mathbb{R}^k \quad \|x_i - x_j\| \sim \delta_{ij}$$

$$X = [x_1 \dots x_n]^T \in \mathbb{R}^{n \times k} \quad d_{ij}(X) = \|x_i - x_j\|$$

$$\Delta \rightsquigarrow D(X) = (\|x_i - x_j\|)_{\substack{i,j \in n \\ i < j}}$$

$$\downarrow$$

$$\Delta^{(q)} = (\delta_{ij}^{(q)})_{i,j} \quad D(X)^{(q)} = (\|x_i - x_j\|^q)_{i,j}$$

Our Goal: $\min_{X \in \mathbb{R}^{n \times k}} \| \Delta^{(q)} - D(X)^{(q)} \|$

$\|\cdot\|$ matrix norm.

4.2.1. Characterization of Euclidean Distance Matrix

$$D^{(ch)}(X) \quad X = [\alpha_1 \dots \alpha_n]^T \in \mathbb{R}^{n \times k} \quad \Delta^{(ch)}$$

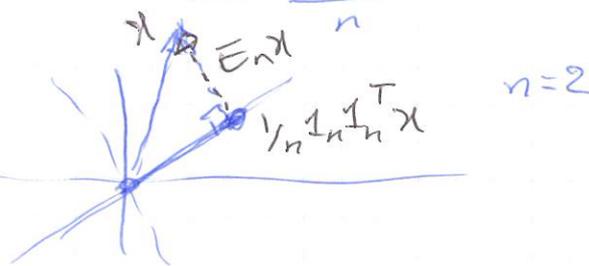
$$E_n = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \quad \mathbf{1}_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

Centering matrix

$$\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T X = \frac{1}{n} \mathbf{1}_n \left(\sum_{i=1}^n \alpha_i \right) = \begin{bmatrix} \bar{\alpha} \\ \vdots \\ \bar{\alpha} \end{bmatrix}$$

$$E_n X = X - \begin{bmatrix} \bar{\alpha} \\ \vdots \\ \bar{\alpha} \end{bmatrix}$$

$$\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$$



Theorem 4.3. The dissimilarity matrix $\Delta \in \mathbb{Z}_n$ has an Euclidean embedding in \mathbb{R}^k if and only

if $\frac{-1}{2} E_n \Delta^{(2)} E_n$ is non-negative definite and

$\text{rk} \left(\frac{-1}{2} E_n \Delta^{(2)} E_n \right) \leq k$. The least k which allows for an embedding is called dimensionality of Δ .

Proof: $X = [x_1 \dots x_n]^T \in \mathbb{R}^{n \times K}$

$$-\frac{1}{2} D^{(2)}(X) = XX^T - \mathbf{1}_n \hat{\lambda}^T - \hat{\lambda} \mathbf{1}_n^T \quad (\text{exercise})$$

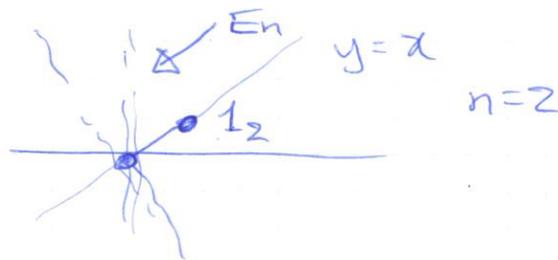
$$\hat{\lambda} = \frac{1}{2} [x_1^T x_1 \dots x_n^T x_n]^T$$

Hint: $\|x_i - x_j\|^2 = (x_i - x_j)^T (x_i - x_j)$
 $= x_i^T x_i + x_j^T x_j - 2x_i^T x_j$

$$-\frac{1}{2} E_n D^{(2)}(X) E_n = E_n X X^T E_n$$

$$E_n \mathbf{1}_n = 0$$

$$\mathbf{1}_n^T E_n = 0$$



first: $-\frac{1}{2} E_n D^{(2)}(X) E_n = E_n X (E_n X)^T \neq 0$

second: $\text{rk}(-\frac{1}{2} E_n D^{(2)}(X) E_n) = \text{rk}(E_n X X^T E_n) \leq K$

$\swarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $n \times n \quad n \times K \quad K \times n \quad n \times n$

if there is X for diss. matrix $\Delta = \begin{matrix} n & \leftarrow & n & \leftarrow & K & \leftarrow & n & \leftarrow & n \\ & & & & \underline{=} & & & & \\ & & & & K & & & & \end{matrix} \subseteq n$

$$\Delta^{(2)} = D^{(2)}(X)$$

$$\Rightarrow -\frac{1}{2} E_n \Delta^{(2)} E_n = -\frac{1}{2} E_n D^{(2)}(X) E_n$$

- ① non.d.
- ② $\text{rk}(-\frac{1}{2} E_n D^{(2)}(X) E_n) \leq K$

$$-\frac{1}{2} E_n \Delta^{(2)} E_n \succeq 0 \quad \text{rk}\left(-\frac{1}{2} E_n \Delta^{(2)} E_n\right) \leq K$$

$$X = ?$$

$$-\frac{1}{2} E_n \Delta^{(2)} E_n = E_n X X^T E_n \\ (E_n X)(E_n X)^T$$

There exists $n \times K$ matrix X such that

$$-\frac{1}{2} E_n \Delta^{(2)} E_n = X X^T \quad X^T E_n = X^T \quad (\text{Proof})$$

Then $\lambda = [\lambda_1 \dots \lambda_n]^T$ is an appropriate

configuration:
$$-\frac{1}{2} E_n \Delta^{(2)} (X) E_n = E_n X X^T E_n = X X^T \\ = -\frac{1}{2} E_n \Delta^{(2)} E_n.$$

Step 1) find $-\frac{1}{2} E_n \Delta^{(2)} E_n$ □

Step 2) See if it is n.n.d. and check the rank.

Step 3) if it is non-negative definite Then

find X :
$$-\frac{1}{2} E_n \Delta^{(2)} E_n = X X^T.$$

Hint: $A \succeq 0 \Rightarrow A = V \Lambda V^T = \underbrace{(V \Lambda^{1/2})}_X (V \Lambda^{1/2} V^T)$

4.2.2. The Best Euclidean Fit to a given dissimilarity matrix.

$$\Delta \in \mathcal{U}_n \quad d \rightarrow \mathcal{D}(X)$$

notation $\lambda^+ = \max(\lambda, 0)$

Theorem 4.4. Let $\Delta \in \mathcal{U}_n$ be a dissimilarity matrix and $\frac{1}{2} E_n \Delta^{(2)} E_n$ has the spectral

decomposition: $\frac{1}{2} E_n \Delta^{(2)} E_n = V \text{diag}(\lambda_1, \dots, \lambda_n) V^T$

with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ and V an orthogonal matrix.

Then the optimization problem

$$\min_{X \in \mathbb{R}^{n \times k}} \| E_n (\Delta^{(2)} - \mathcal{D}^{(2)}(X)) E_n \|$$

has a solution given by

$$X^* = [\sqrt{\lambda_1^+} v_1, \dots, \sqrt{\lambda_k^+} v_k] \in \mathbb{R}^{n \times k}$$

Proof. Note that a solution to

$$\min_{\substack{A \succeq 0 \\ \text{rk}(A) \leq k}} \| \frac{1}{2} E_n \Delta^{(2)} E_n - A \|$$

(Thm. 2.6)

is given by

$$A^* = V \text{diag}(\lambda_1^+, \dots, \lambda_k^+, 0, \dots, 0) V^T.$$

$$\begin{aligned} -\frac{1}{2} E_n D^{(2)}(X^*) E_n &= E_n X^* X^{*T} E_n \\ &= E_n [v_1 \dots v_k] \text{diag}(\lambda_1^+, \dots, \lambda_k^+) \\ &\quad X [v_1 \dots v_k]^T E_n \\ &= V \text{diag}(\lambda_1^+, \dots, \lambda_k^+, 0, \dots, 0) V^T = A^* \end{aligned}$$

$$\Rightarrow X^* \rightarrow D^{(2)}(X^*) \rightarrow -\frac{1}{2} E_n D^{(2)}(X^*) E_n = A^*$$

$$X^* \rightarrow \underset{\text{minimizer}}{\| -\frac{1}{2} E_n (\Delta^{(2)} - D^{(2)}(X)) E_n \|} \quad 0$$

Step 1) $-\frac{1}{2} E_n \Delta^{(2)} E_n$

Step 2) Spectral decomposition

Step 3) $[\sqrt{\lambda_1^+} v_1, \dots, \sqrt{\lambda_k^+} v_k]$

4.2.3 non-linear Dimensionality Reduction

- Manifold learning
- Geodesic

Isomap: complete isometric feature mapping
(Zou)

$$x_1, \dots, x_n \in \mathbb{R}^p$$

1) Construct a neighborhood graph: find a weighted graph $G(V, E, W)$

$v_i = x_i$. Two vertices v_i and v_j are connected only if $\|x_i - x_j\| \leq \epsilon$.

(another way is to connect each point to its k -nearest neighbors)

2) Compute the shortest paths for each pair (v_i, v_j) (Dijkstra's algorithm).

The geodesic distance $\delta_{ij} = \delta(v_i, v_j)$ can be taken as number of hops/links from v_i to v_j or sum of $\|x_k - x_l\|$ on the shortest path.

3) Construct d -dimensional embedding -

Apply MDS on the basis of

geodesic distances $\Delta = (\delta_{ij})_{1 \leq i, j \leq n}$.

1) Very large distance may distort local neighborhoods

2) Computational Complexity

3) Not robust to noise : choice of " ϵ "

4.3) Diffusion maps.

Weights of $G(V, E, W)$

weight functions or kernel $w_{ij} = K(x_i, x_j)$

The selected Kernel should satisfy three properties

- Symmetry $K(x_i, x_j) = K(x_j, x_i)$

- non-negative $K(x_i, x_j) \geq 0$

- Locality : there is a scale parameter ϵ

such that if $\|x_i - x_j\| \ll \epsilon$ then $K(x_i, x_j) \rightarrow 1$

if $\|x_i - x_j\| \gg \epsilon$ then $K(x_i, x_j) \rightarrow 0$

Gaussian Kernel:

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\epsilon^2}\right)$$

Caveat: Do not confuse this "Kernel"

with kernels of support vector machines
(later!)