

# Clustering Matrices: A Metric Learning Approach to Disease Subtyping in Mental Health

Hanchao Zhang

Division of Biostatistics, Department of Population Health,  
Grossman School of Medicine, New York University

March 21, 2023

## Acknowledgements

### Collaborators

This talk presents joint work with the following:

- ▶ Thaddeus Tarpey (Grossman School of Medicine, NYU)
- ▶ Emily R. Stern (Grossman School of Medicine, NYU)
- ▶ Alessandro S. De Nadai (Harvard Medical School)

### Fundings

This work is supported by NIMH grants:

- ▶ R01 MH099003
- ▶ R01 MH126981
- ▶ R01 MH111794
- ▶ R01 MH111794
- ▶ R61/R33 MH107589

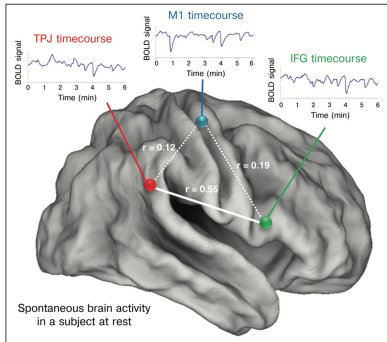
## Outline for section 1

- 1 Introduction
- 2 Distance Metrics for Clustering Matrices
- 3 Models for Clustering Matrices
  - Unconstrained Common Principal Components
  - Partial Common Principal Components
- 4 Simulation
- 5 Discussion

# Functional Connectivity Matrix

## Functional Connectivity

Functional connectivity is defined as the temporal coincidence of spatially distant neurophysiological events.



Functional Connectivity (Gillebert and Mantini, 2013)

## Functional Connectivity Matrix

For each participant  $i$ , let  $y_{ij} \in \mathbb{R}^T$  be the longitudinal measurement of blood oxygen level-dependent (BOLD) signal on the region of interest  $j$ ,  $j = 1, 2, \dots, p$ .

The functional connectivity matrix for participant  $i$ :  $\Sigma_i = \text{Cov}(y_i) \succeq \mathbf{0}$

## Self-Consistency Clustering Algorithms

### Scalar Outcomes

K-Means algorithm (Steinhaus et al., 1956)

$$\text{minimize}_C \underbrace{\frac{1}{n} \sum_{i=1}^m \sum_{k \in C_i} \|x_k - \bar{x}_{C_i}\|^2}_{\text{within cluster sum of squares}} \quad \text{OR} \quad \text{maximize}_C \underbrace{\sum_{i=1}^m \frac{n_{C_i}}{n} \cdot \|\bar{x}_{C_i}\|^2}_{\text{between cluster sum of squares}}$$

### Functional Outcomes

Clustering Functional data (Tarpey and Kinaterder, 2003)

$y_i(t), i = 1, \dots, n, t \in T$ , typically a compact real interval,  $y_i(t) =$  function

$$y_i(t) = \mathbf{b}'(t) \beta_i + \epsilon_i(t) = \sum_{j=1}^{\infty} \beta_{ij} b_j(t) + \epsilon_i(t).$$

- ▶  $\mathbf{b} = (b_1(t), \dots, b_p(t), \dots)'$  is basis functions
- ▶  $\beta_i = (\beta_{1i}, \dots, \beta_{ip}, \dots)'$  is a vector of basis coefficients

Perform K-Means or other algorithms on basis coefficients  $\beta_i$

## Positive Semi-Definite Matrix Outcomes

Consider that for each observation  $i$ ,  $i = 1, 2, \dots, n$ , we observe  $p$  functional outcomes  $y_{ij}(t)$  with mean 0,  $j = 1, 2, \dots, p$ . Then we can obtain a positive semi-definite matrix for subject  $i$ :

$$\Psi_i = \int_T \mathbf{y}_i(t)^T \mathbf{y}_i(t) dt, \Psi_i \succcurlyeq 0,$$

where  $\Psi_i \succcurlyeq 0$  means  $\Psi_i$  is positive semi-definite matrix. (All the eigenvalues of  $\Psi_i$  are larger and equal to 0).

## Clustering Algorithm Approaches

- ▶ cluster subjects by  $\Psi_i$ 's,  $i = 1, 2, \dots, n$
- ▶ vectorize  $\Psi_i$ 's and treat it as vector
- ▶ consider some distance metrics for matrix similarity
- ▶ consider the probability distribution (e.g., Wishart distribution)

## Outline for section 2

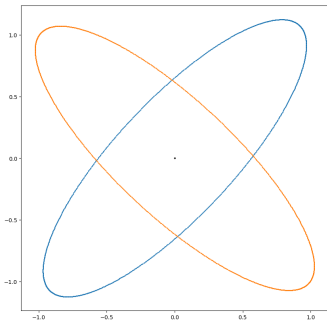
- 1 Introduction
- 2 Distance Metrics for Clustering Matrices
- 3 Models for Clustering Matrices
  - Unconstrained Common Principal Components
  - Partial Common Principal Components
- 4 Simulation
- 5 Discussion

## Distance Metrics for Matrices

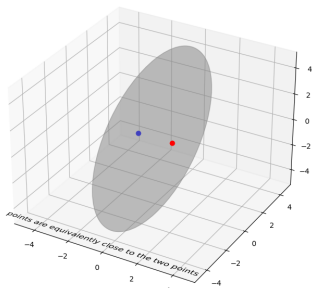
### Euclidean Distance (chapter 2 Minh and Murino, 2017)

let  $\Psi_1$  and  $\Psi_2 \in \mathbb{R}^{2 \times 2}$  be two positive semi-definite matrices. The Euclidean distance between two matrices  $d_E(\Psi_1, \Psi_2)$  can be represented by points in  $\mathbb{R}^3$

$$d_E(\Psi_1, \Psi_2) = \|\Psi_1 - \Psi_2\|_F^2 = \|\text{vec}(\Psi_1^T) - \text{vec}(\Psi_2^T)\|^2$$



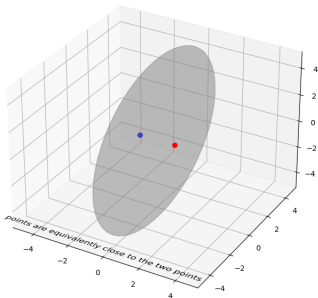
2 positive semi-definite matrices



the vectorized matrices



## Distance Metrics for Matrices



the vectorized matrices

### Disadvantage of Euclidean Distance

- ▶ matrices with similar shapes are clustered into different groups

## Distance Metrics for Matrices (chapter 2 Malhi et al., 2017)

### Other Metrics

### Affine-invariant Riemannian Distance

$$d_{aiE}(\mathbf{A}, \mathbf{B}) = \|\log(\mathbf{A}^{-\frac{1}{2}} \mathbf{B} \mathbf{A}^{-\frac{1}{2}})\|_F$$

### Log-Determinant Divergences

$$d_{\log \det}^1(\mathbf{A}, \mathbf{B}) = \text{tr}(\mathbf{B}^{-1} \mathbf{A} - \mathbf{I}) - \log \det(\mathbf{B}^{-1} \mathbf{A})$$

### Symmetric Stein Divergence

$$d_{\text{stein}}^2(\mathbf{A}, \mathbf{B}) = \log \det\left(\frac{\mathbf{A} + \mathbf{B}}{2}\right) - \frac{1}{2} \log \det(\mathbf{A} \mathbf{B})$$

### Disadvantages

- ▶ does not consider the structure (shape) of p.s.d matrices

## Outline for section 3

- 1 Introduction
- 2 Distance Metrics for Clustering Matrices
- 3 Models for Clustering Matrices**
  - Unconstrained Common Principal Components
  - Partial Common Principal Components
- 4 Simulation
- 5 Discussion

## Common Principal Components (CPC) (Flury, 1984)

### Definition

Let  $\Psi_1, \dots, \Psi_n$  be positive definite symmetric matrix of dimension  $p \times p$ , we wish to find an orthonormal matrix  $B$  which makes the  $\Psi_i$ 's simultaneously "as diagonal as possible":

### Objective Function

Let  $F_i$  be the transformed  $\Psi_i$  by  $B$ :

$$F_i = B^T \Psi_i B$$

To make sure that  $F_i$ 's,  $i = 1, \dots, n$  are as diagonal as possible, we wish to minimize:

$$\underset{B}{\text{minimize}} \prod_{i=1}^n \left\{ \frac{\det(\text{diag}(F_i))}{\det(F_i)} \right\} = \prod_{i=1}^n \left\{ \frac{\det(\text{diag}(B^T \Psi_i B))}{\det(B^T \Psi_i B)} \right\},$$

where  $\det(F_i) \leq \det(\text{diag}(F_i))$ .

### Algorithms

- ▶ FG-algorithm (Flury and Constantine, 1985)
- ▶ MM algorithms (Browne and McNicholas, 2014)
- ▶ R algorithm (Hallin et al., 2014)

# Unconstrained Common Principal Components for Clustering Matrices

## Stiefel Manifold

the Stiefel manifold  $V_k(\mathbb{R}^P)$  is the set of all orthonormal  $k$ -frames in  $\mathbb{R}^P$

$$V_k(\mathbb{R}^P) = \{\mathbf{A} \in \mathbb{R}^{n \times P} : \mathbf{A}^T \mathbf{A} = \mathbf{I}_P\}$$

## Self-Consistency Algorithm Based on CPC

Let  $S \subset \mathbb{R}^{P \times P}$ ,  $S \succeq 0$  denote a set of p.s.d. matrices. For each  $\mathbf{B} \in V_p(\mathbb{R}^P)$ , define:

$$D_{\mathbf{B}}(S) = \{\Psi \in \mathbb{R}^{P \times P}, \Psi \succeq 0 : \|\Psi - \underbrace{\mathbf{B} \text{diag}(\mathbf{F}_{\mathbf{B}}) \mathbf{B}^T}_{\hat{\Psi}}\|_F \leq \|\Psi - \mathbf{A} \text{diag}(\mathbf{F}_{\mathbf{A}}) \mathbf{A}^T\|_F, \\ \mathbf{B} \neq \mathbf{A}, \mathbf{A} \in V_p(\mathbb{R}^P)\}.$$

Therefore, each matrix in set  $D_{\mathbf{B}}(S)$  shares common principal components  $\mathbf{B}$  that can make them “as diagonal as possible”.

## Note

Since we have  $\Psi = \mathbf{B}\mathbf{F}_{\mathbf{B}}\mathbf{B}^T = \mathbf{A}\mathbf{F}_{\mathbf{A}}\mathbf{A}^T$ , we can redefine  $D_{\mathbf{B}}(S)$  as follow:

$$D_{\mathbf{B}}(S) = \{\Psi \in \mathbb{R}^{P \times P}, \Psi \succeq 0 : \|\mathbf{F}_{\mathbf{B}} - \text{diag}(\mathbf{F}_{\mathbf{B}})\|_F \leq \|\mathbf{F}_{\mathbf{A}} - \text{diag}(\mathbf{F}_{\mathbf{A}})\|_F, \\ \mathbf{B} \neq \mathbf{A}, \mathbf{A} \in V_p(\mathbb{R}^P)\}.$$

# Unconstrained Common Principal Components for Clustering Matrices

## Self-Consistency Algorithm Based on CPC

Let  $S \in \mathbb{R}^{p \times p}$ ,  $S \succeq 0$  denote a measurable set. For each  $B \in V_p(\mathbb{R}^p)$ , define:

$$D_B(S) = \{ \Psi \in \mathbb{R}^{p \times p}, \Psi \succeq 0 : \|\Psi - B \operatorname{diag}(F_B) B^T\|_F \leq \|\Psi - A \operatorname{diag}(F_A) A^T\|_F, \\ B \neq A, A \in V_p(\mathbb{R}^p) \}.$$

## Unconstrained CPC for Matrices Clustering

### Algorithm Clustering Matrices Using Unconstrained CPC

Start with an initial partition of all matrices into  $K$  clusters

- 1: for each cluster  $k$ ,  $k = 1, 2, \dots, K$ , estimate the common principal component  $B_k$ .
- 2: assign individual matrices  $\Psi_i$  to cluster  $k$  if

$$k^* = \arg \min_{k=1, \dots, K} \|\Psi_i - B_k \operatorname{diag}(F_{B_k}) B_k^T\|_F$$

repeat steps 1 and 2 until convergence.

## Partial Common Principal Components

### Self-Consistency Algorithm Based on Partial CPC

Let  $S \in \mathbb{R}^{p \times p}$ ,  $S \succeq 0$  denote a measurable set. For each  $\mathbf{B} := (\beta_1, \dots, \beta_m) \in V_m(\mathbb{R}^p)$ ,

$$D_{\mathbf{B}}(S) = \left\{ \Psi \in \mathbb{R}^{p \times p}, \Psi \succeq 0 : \left\| \Psi - \sum_{r=1}^m f_{\mathbf{B}_r} \beta_r \beta_r^T \right\|_F \leq \left\| \Psi - f_{\mathbf{A}_r} \alpha_r \alpha_r^T \right\|_F, \right. \\ \left. \mathbf{B} \neq \mathbf{A}, \mathbf{A} \in V_m(\mathbb{R}^p) \right\},$$

where  $f_1, \dots, f_p$  are the diagonal elements of  $\mathbf{F}$ , and  $m \leq p$ .

### Unconstrained CPC for Matrices Clustering

#### Algorithm Clustering Matrices Using Unconstrained CPC

Start with an initial partition of all matrices into  $K$  clusters

- 1: for each cluster  $k$ ,  $k = 1, 2, \dots, K$ , estimate the common principal component  $\mathbf{B}_k$ .
  - 2: assign individual matrices  $\Psi_i$  to cluster  $k$  if  $k^* = \arg \min_{k=1, \dots, K} \left\| \Psi_i - \sum_{r=1}^m f_{\mathbf{B}_k} \beta_r \beta_r^T \right\|_F$
- repeat steps 1 and 2 until convergence.

## Outline for section 4

- 1 Introduction
- 2 Distance Metrics for Clustering Matrices
- 3 Models for Clustering Matrices
  - Unconstrained Common Principal Components
  - Partial Common Principal Components
- 4 Simulation**
- 5 Discussion



## Simulation

### Simulation Settings

let  $\mathbf{B}_1$  and  $\mathbf{B}_2$  be the two common eigenvectors for the two clusters

$$\mathbf{B}_1 = \begin{pmatrix} \cos(\beta_1) & -\sin(\beta_1) \\ \sin(\beta_1) & \cos(\beta_1) \end{pmatrix} \quad \mathbf{B}_2 = \begin{pmatrix} \cos(\beta_2) & -\sin(\beta_2) \\ \sin(\beta_2) & \cos(\beta_2) \end{pmatrix},$$

where  $\beta_1$  and  $\beta_2$  be 2 scalars from 0 to  $2\pi$ . Let  $|\beta_1 - \beta_2| = \theta$  be the differences between two eigenvectors.

let  $\lambda_{\mathbf{B}_1 i} = [\lambda_{\mathbf{B}_1 i1}, \lambda_{\mathbf{B}_1 i2}]$  and  $\lambda_{\mathbf{B}_2 i} = [\lambda_{\mathbf{B}_2 i1}, \lambda_{\mathbf{B}_2 i2}]$  be the eigenvalues for the two clusters. Denote  $\Lambda_{\mathbf{B}_1 i} = \text{diag}(\lambda_{\mathbf{B}_1 i})$ , and  $\Lambda_{\mathbf{B}_2 i} = \text{diag}(\lambda_{\mathbf{B}_2 i})$ , where  $\lambda \sim \chi^2(\text{df})$ .

Then we can obtain our simulated matrices:

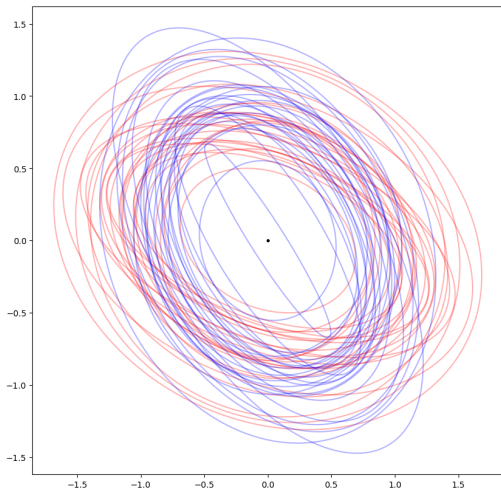
$$\Psi_{1i} = \mathbf{B}_1 \Lambda_{\mathbf{B}_1 i} \mathbf{B}_1^T + \mathbf{E}_1 \quad \Psi_{2i} = \mathbf{B}_2 \Lambda_{\mathbf{B}_2 i} \mathbf{B}_2^T + \mathbf{E}_2,$$

where  $\mathbf{E}_1$ , and  $\mathbf{E}_2$  are random error with mean 0.

### Note

- ▶  $\theta$  denotes how close the two clusters of matrices are
- ▶  $\mathbf{E}$  is some random perturbation on eigenvector and eigenvalues.

# Simulations



$\theta = \pi/5, df = 5$

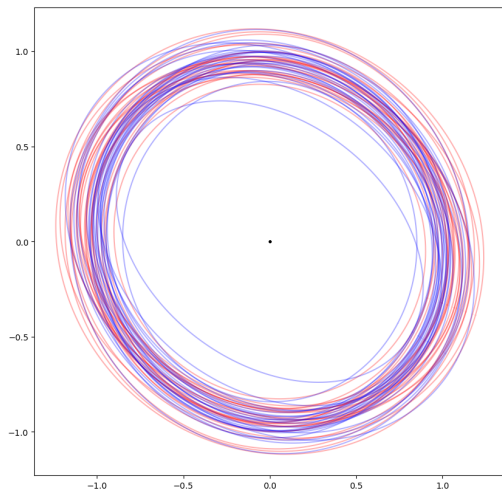
## Simulation

- ▶  $\theta = \pi/5$
- ▶ eigenvalues  $\sim \chi^2(5)$

## Classification Error

- ▶ CPCA = 0
- ▶ rCPCA = 0
- ▶ Euclidean Distance = 0.28
- ▶ Affine Invariance Divergence = 0.34
- ▶ Log-Determinant Divergence = 0.38

## Simulations

 $\theta = \pi/5, df = 5$ 

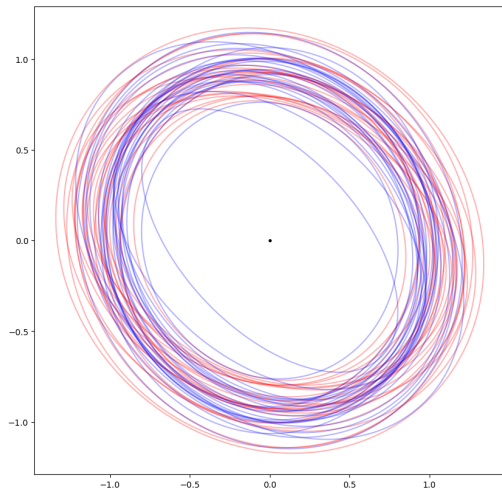
## Simulation

- ▶  $\theta = \pi/15$
- ▶ eigenvalues  $\sim \chi^2(40)$

## Classification Error

- ▶ CPCA = 0
- ▶ rCPCA = 0.02
- ▶ Euclidean Distance = 0.4
- ▶ Affine Invariance Divergence = 0.4
- ▶ Log-Determinant Divergence = 0.44

# Simulations



$\theta = \pi/5, df = 5$

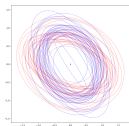
## Simulation

- ▶  $\theta = \pi/15$
- ▶ eigenvalues  $\sim \chi^2(20)$
- ▶ noise = 5%

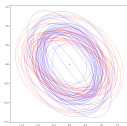
## Classification Error

- ▶ CPCA = 0.02
- ▶ rCPCA = 0.02
- ▶ Euclidean Distance = 0.4
- ▶ Affine Invariance Divergence = 0.4
- ▶ Log-Determinant Divergence = 0.44

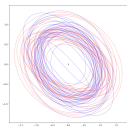
# Simulations



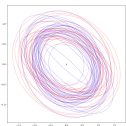
$$\theta = \pi/5, df = 5$$



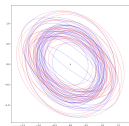
$$\theta = \pi/6, df = 5$$



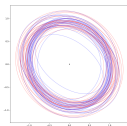
$$\theta = \pi/8, df = 5$$



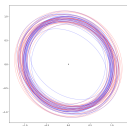
$$\theta = \pi/10, df = 5$$



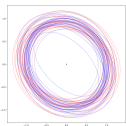
$$\theta = \pi/15, df = 5$$



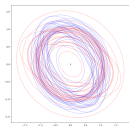
$$\theta = \pi/15, df = 20$$



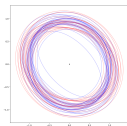
$$\theta = \pi/15, df = 40$$



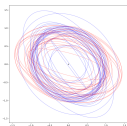
$$\theta = \pi/15, df = 20, \text{ noise} = 1\%$$



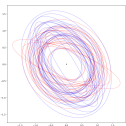
$$\theta = \pi/15, df = 20, \text{ noise} = 2\%$$



$$\theta = \pi/15, df = 20, \text{ noise} = 4\%$$



$$\theta = \pi/15, df = 20, \text{ noise} = 6\%$$



$$\theta = \pi/15, df = 20, \text{ noise} = 10\%$$

# Simulations

## Simulation Results

Classification Errors					
Methods	$\theta = \pi/5$ $df = 5$	$\theta = \pi/6$ $df = 5$	$\theta = \pi/8$ $df = 5$	$\theta = \pi/10$ $df = 5$	$\theta = \pi/15$ $df = 5$
CPCA	0.00	0.00	0.00	0.00	0.00
rCPCA	0.00	0.00	0.02	0.00	0.02
Frobenius	0.28	0.32	0.32	0.34	0.34
Aff. Div.	0.34	0.38	0.42	0.42	0.44
Log-Det	0.38	0.38	0.44	0.44	0.46
Methods	$\theta = \pi/15$ $df = 20$	$\theta = \pi/15$ $df = 40$	$\theta = \pi/15$ $df = 20$ noise = 0.01	$\theta = \pi/15$ $df = 20$ noise = 0.06	$\theta = \pi/15$ $df = 20$ noise = 0.10
CPCA	0.00	0.00	0.02	0.10	0.18
rCPCA	0.10	0.12	0.02	0.14	0.20
Frobenius	0.40	0.40	0.40	0.36	0.50
Aff. Div.	0.40	0.40	0.40	0.40	0.46
Log-Det	0.44	0.44	0.44	0.44	0.48

## Outline for section 5

- 1 Introduction
- 2 Distance Metrics for Clustering Matrices
- 3 Models for Clustering Matrices
  - Unconstrained Common Principal Components
  - Partial Common Principal Components
- 4 Simulation
- 5 Discussion

## Referenes

- R. P. Browne and P. D. McNicholas. Estimating common principal components in high dimensions. *Advances in Data Analysis and Classification*, 8:217–226, 2014.
- B. N. Flury. Common principal components in k groups. *Journal of the American Statistical Association*, 79(388):892–898, 1984.
- B. N. Flury and G. Constantine. Algorithm as 211: The fg diagonalization algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 34(2):177–183, 1985.
- C. R. Gillebert and D. Mantini. Functional connectivity in the normal and injured brain. *The Neuroscientist*, 19(5):509–522, 2013.
- M. Hallin, D. Paindaveine, and T. Verdebout. Efficient r-estimation of principal and common principal components. *Journal of the American Statistical Association*, 109(507):1071–1083, 2014.
- G. S. Malhi, Y. Byrow, T. Outhred, and K. Fritz. Exclusion of overlapping symptoms in dsm-5 mixed features specifier: heuristic diagnostic and treatment implications. *CNS spectrums*, 22(2):126–133, 2017.
- H. Q. Minh and V. Murino. Covariances in computer vision and machine learning. *Synthesis Lectures on Computer Vision*, 7(4):1–170, 2017.
- H. Steinhaus et al. Sur la division des corps matériels en parties. *Bull. Acad. Polon. Sci*, 1(804):801, 1956.
- T. Tarpey and K. K. Kinatader. Clustering functional data. *Journal of classification*, 20(1), 2003.