Structured data analysis with RGCCA

Arthur Tenenhaus

Séminaire de Statistique appliquée du CNAM, 16/03/2018













Spinocerebellar Ataxia (SCA) dataset

Brain and Spine Institute



The objective of multiblock component methods is to find block components summarizing the relevant information between and within the blocks.

Volume of the pons vs groups



Objective. Identify a set of variables within each block that influence the volume of the pons while accounting for the multiblock structure of the SCA dataset.

The SCA dataset from an RGCCA point of view



The philosophy of multiblock component methods



The philosophy of multiblock component methods



Block components should verified two properties at the same time:

- 1. Block components well explain their own block.
- 2. Block components are as correlated as possible for connected blocks.

The philosophy of multiblock component methods



BLOCKS ARE PARTIALLY CONNECTED $c_{jk} = 1 \text{ if } \mathbf{X}_j \leftrightarrow \mathbf{X}_k, 0 \text{ otherwise}$ $\mathbf{X}_j = \mathbf{X}_j$		
SUMCOR	$\max_{\operatorname{var}(\mathbf{X}_{j}\mathbf{w}_{j})=1} \sum_{j,k} c_{jk} \operatorname{cov}(\mathbf{X}_{j}\mathbf{w}_{j}, \mathbf{X}_{k}\mathbf{w}_{k})$	
SSQCOR	$\max_{\operatorname{var}(\mathbf{X}_{j}\mathbf{w}_{j})=1} \sum_{j,k} c_{jk} \operatorname{cov}^{2}(\mathbf{X}_{j}\mathbf{w}_{j}, \mathbf{X}_{k}\mathbf{w}_{k})$	
SABSCOR	$\max_{\operatorname{var}(\mathbf{X}_{j}\mathbf{w}_{j})=1}\sum_{j,k}\frac{c_{jk}}{ \operatorname{cov}(\mathbf{X}_{j}\mathbf{w}_{j},\mathbf{X}_{k}\mathbf{w}_{k}) }$	

RGCCA for multiblock analysis

$$\max_{\mathbf{w}_{1},...,\mathbf{w}_{J}} \sum_{j,k}^{J} c_{jk} g\left(\operatorname{cov}(\mathbf{X}_{j} \mathbf{w}_{j}, \mathbf{X}_{k} \mathbf{w}_{k}) \right)$$

s.t. $(1 - \tau_{j}) \operatorname{var}(\mathbf{X}_{j} \mathbf{w}_{j}) + \tau_{j} \left\| \mathbf{w}_{j} \right\|_{2}^{2} = 1, j = 1, ..., J$

• $c_{jk} = 1$ if $\mathbf{X}_j \leftrightarrow \mathbf{X}_k$, 0 otherwise

•
$$g =$$
 any convex function – e.g.
$$\begin{cases} g(x) = x & (\text{Horst sheme}) \\ g(x) = x^2 & (\text{Factorial scheme}) \\ g(x) = |x| & (\text{Centroid scheme}) \end{cases}$$

• $0 \le \tau_j \le 1$ continuum between correlation and covariance

Tenenhaus A. and Guillemot V. (2017): RGCCA Package. http://cran.project.org/web/packages/RGCCA/index.html

Tenenhaus M, Tenenhaus A, Groenen PJF, (2017) Regularized generalized canonical correlation analysis: A framework for sequential multiblock component methods, Psychometrika, vol. 82, no. 3, 737–777

Tenenhaus A, Philippe C, Frouin V (2015) Kernel generalized canonical correlation analysis, Computational Statistics & Data Analysis, vol. 90, pp. 114-131. Tenenhaus A, Tenenhaus M (2011) Regularized generalized canonical correlation analysis, vol. 76, pp. 257-284, Psychometrika.

BLOCKS ARE CONNECTED TO THE SUPERBLOCK $X_{J+1} = [X_1, ..., X_J]$

 \mathbf{X}_1

 \mathbf{X}_{I+1}



2-block special cases

$$\max_{\mathbf{w}_1,\mathbf{w}_2} \operatorname{cov}(\mathbf{X}_1\mathbf{w}_1,\mathbf{X}_2\mathbf{w}_2) \quad \text{s.t.} \quad (1-\tau_j)\operatorname{var}(\mathbf{X}_j\mathbf{w}_j) + \tau_j \|\mathbf{w}_j\|_2^2 = 1, j = 1, 2$$

Methods	Criterion	Constraints
PLS regression	$\max \operatorname{cov}(\mathbf{X}_1\mathbf{w}_1, \mathbf{X}_2\mathbf{w}_2)$	$\ \mathbf{w}_1\ = \ \mathbf{w}_2\ = 1$
Canonical Correlation Analysis	$\max \operatorname{cor}(\mathbf{X}_1 \mathbf{w}_1, \mathbf{X}_2 \mathbf{w}_2)$	$var(\mathbf{X}_1\mathbf{w}_1) = var(\mathbf{X}_2\mathbf{w}_2) = 1$
Redundancy analysis of X1 with respect to X2	$\max \operatorname{cor}(\mathbf{X}_1 \mathbf{w}_1, \mathbf{X}_2 \mathbf{w}_2) \times \operatorname{var}(\mathbf{X}_1 \mathbf{w}_1)$	$var(X_2w_2) = 1$ $\ w_1\ = 1$

Components $X_1 w_1$ and $X_2 w_2$ $X_2 a_2$ are well correlated. 1st component is stable No stability condition for 2nd component

Hierarchical PCA with RGCCA

Hierarchical PCA optimization problem:

$$\max_{\substack{\|\mathbf{a}_{j}\|=\cdots\|\mathbf{a}_{J}\|=1\\ \operatorname{var}(X_{J+1}\mathbf{a}_{J+1})=1}} \sum_{j=1}^{J} \operatorname{cov}^{4}(\mathbf{X}_{j}\mathbf{w}_{j}, \mathbf{X}_{J+1}\mathbf{w}_{J+1})$$

• Scheme function: $g(x) = x^4$

• Shrinkage parameters: $\tau_1 = \cdots = \tau_J = 1$ and $\tau_{J+1} = 0$



Choice of the shrinkage constant τ_i

$$\max_{\mathbf{a}_{1},...,\mathbf{a}_{J}} \sum_{j,k}^{J} c_{jk} g\left(\operatorname{cov}(\mathbf{X}_{j} \mathbf{w}_{j}, \mathbf{X}_{k} \mathbf{w}_{k})\right)$$

s.t. $\mathbf{w}_{j}^{\mathsf{T}} \left((1 - \tau_{j}) n^{-1} \mathbf{X}_{j}^{\mathsf{T}} \mathbf{X}_{j} + \tau_{j} \mathbf{I}_{p_{j}}\right) \mathbf{w}_{j} = 1, j = 1, ..., J$



RGCCA for multiblock analysis

$$\max_{\mathbf{w}_{1},...,\mathbf{w}_{J}} \sum_{j,k}^{J} c_{jk} g\left(\operatorname{cov}(\mathbf{X}_{j} \mathbf{w}_{j}, \mathbf{X}_{k} \mathbf{w}_{k})\right)$$

s.t. $(1 - \tau_{j})\operatorname{var}(\mathbf{X}_{j} \mathbf{w}_{j}) + \tau_{j} \|\mathbf{w}_{j}\|_{2}^{2} = 1, j = 1, ..., J$

Two key ingredients:

(i) Block relaxation

(ii) Majorization by Minorization (MM)

Tenenhaus A. and Guillemot V. (2017): RGCCA Package. http://cran.project.org/web/packages/RGCCA/index.html

Tenenhaus M, Tenenhaus A, Groenen PJF, (2017) Regularized generalized canonical correlation analysis: A framework for sequential multiblock component methods, Psychometrika, vol. 82, no. 3, 737–777

Tenenhaus A, Philippe C, Frouin V (2015) Kernel generalized canonical correlation analysis, Computational Statistics & Data Analysis, vol. 90, pp. 114-131. Tenenhaus A, Tenenhaus M (2011) Regularized generalized canonical correlation analysis, vol. 76, pp. 257-284, Psychometrika.

Block relaxation: from w^s to w^{s+1}

$$\boldsymbol{w}^{s} = \left(\mathbf{w}_{1}^{s}, \mathbf{w}_{2}^{s}, \dots, \mathbf{w}_{J}^{s} \right)$$

$$\underset{\mathbf{w}_{1},\mathbf{w}_{1}^{\mathsf{T}}\mathsf{M}\mathbf{w}_{1}=1}{\operatorname{argmax}} h(\mathbf{w}_{1},\mathbf{w}_{2}^{s},\ldots,\mathbf{w}_{J}^{s}) \longrightarrow \mathbf{w}_{1}^{s+1}$$

.

 $\rightarrow \mathbf{w}_{j}^{s+1}$

 $\rightarrow \mathbf{W}_{J}^{s+1}$

MM approach for RGCCA



• Tenenhaus A. and Tenenhaus M., Regularized Generalized Canonical Correlation Analysis, Psychometrika, vol. 76, Issue 2, pp. 257-284, 2011

• Tenenhaus M, Tenenhaus A, Groenen P.J.F, (2017) Regularized generalized canonical correlation analysis: A framework for sequential multiblock component methods, Psychometrika, doi: 10.1007/s11336-017-9573-x

• Tenenhaus, A., Philippe, C., Frouin, V. (2015). Kernel Generalized Canonical Correlation Analysis. Computational Statistics & Data Analysis, 90, 114-131.

• Tenenhaus A. and Guillemot V. (2017): RGCCA Package. http://cran.project.org/web/packages/RGCCA/index.html

Properties of the RGCCA algorithm for multiblock data

► Monotone convergence: $h(\mathbf{w}_1^{s+1}, ..., \mathbf{w}_J^{s+1}) \ge h(\mathbf{w}_1^s, ..., \mathbf{w}_J^s)$.

In addition, assuming uniqueness of the solution of the MM step, the following properties hold:

- The sequence $\{\mathbf{w}^s\}$ is asymptotically regular: $\lim_{s\to\infty} ||\mathbf{w}^{s+1} \mathbf{w}^s|| = 0$.
- ► At convergence, a stationary point is obtained.

[•] Tenenhaus M, Tenenhaus A, Groenen P.J.F, (2017) Regularized generalized canonical correlation analysis: A framework for sequential multiblock component methods, Psychometrika, doi: 10.1007/s11336-017-9573-x

[•] Tenenhaus A. and Guillemot V. (2017): RGCCA Package. <u>http://cran.project.org/web/packages/RGCCA/index.html</u>

multiblock component methods with sparsity



Block components should verified two properties at the same time:

- 1. Block components well explain their own block.
- 2. Block components are as correlated as possible for connected blocks.

3. Block components are built from sparse a_j

extension of RGCCA for sparse solution

$$\max_{\mathbf{w}_{j}} h(\mathbf{w}_{1}^{s+1}, \dots, \mathbf{w}_{j-1}^{s+1}, \mathbf{w}_{j}, \mathbf{w}_{j+1}^{s}, \dots, \mathbf{w}_{j}^{s}) \quad \text{s.t.} \quad \left\|\mathbf{w}_{j}\right\|_{2}^{2} = 1 \quad \& \quad \left\|\mathbf{w}_{j}\right\|_{1} \leq s_{j}$$
sparsity constant
$$h(\mathbf{w}_{1}^{s+1}, \dots, \mathbf{w}_{j-1}^{s+1}, \mathbf{w}_{j}, \mathbf{w}_{j+1}^{s}, \dots, \mathbf{w}_{j}^{s})$$

$$\tilde{f}_{j}(\mathbf{w}_{j}, \mathbf{w}_{-j}^{s}) \quad \text{s.t.} \quad \left\|\mathbf{w}_{j}\right\|_{2} = 1 \text{ and } \left\|\mathbf{w}_{j}\right\|_{1} \leq s_{j}$$
obtained by soft-thresholding
$$(\mathbf{w}_{j}^{s+1} - \mathbf{w}_{j}^{s})$$

Tenenhaus A., Philippe C., Guillemot V, et al..(2014). Variable Selection for Generalized Canonical Correlation Analysis, Biostatistics, 15 (3): 569-583 Tenenhaus A. and Guillemot V. (2017): RGCCA Package. <u>http://cran.project.org/web/packages/RGCCA/index.html</u>

The SCA dataset from an RGCCA point of view



Divide and conquer strategy with SGCCA

Arginine_Proline



Visualization















The COMA project from a multiblock viewpoint

(Brain and Spine Institute, La pitié Salpêtrière Hospital)

Predict the long term recovery of patients after traumatic brain injury



The COMA project from a multiblock viewpoint

(Brain and Spine Institute, La pitié Salpêtrière Hospital)

Predict the long term recovery of patients after traumatic brain injury



... to Multiblock / Multiway data



MGCCA optimization problem



The MGCCA algorithm

$$\max_{\mathbf{w}_1,\dots,\mathbf{w}_J} \sum_{j,k=1}^J c_{jk} g\left(\operatorname{cov}(\mathbf{X}_j \mathbf{w}_j, \mathbf{X}_k \mathbf{w}_k) \right) \quad \text{s.t.} \quad \mathbf{w}_j^\top \mathbf{M}_j \mathbf{w}_j = 1 \text{ and } \mathbf{w}_j = \mathbf{w}_j^K \otimes \mathbf{w}_j^J, j = 1, \dots, J$$

Two key ingredients: Block relaxation and Majorization by Minorization



• Tenenhaus, A., Le Brusquet, L., Lechuga, G. (2015). Multiway Regularized Generalized Canonical Correlation Analysis. 47ème Journées de Statistique, Lille, France

• Tenenhaus A., Le Brusquet L. Three-way Regularized Generalized Canonical Correlation Analysis, ThRee-way methods In Chemistry And Psychology, (TRICAP), 2015

MGCCA results

Predict the long term recovery of patients after traumatic brain injury



Influence of spatial positions



Discriminating voxels within the white matter bundles

Influence of spatial positions

Modality	\mathbf{w}_1^K
FA	09887
MD	0,0036
L ₁	0,0046
L _t	0,0031







RGCCA for multiblock data analysis

Block components should verify two properties at the same time:

- (i) Block components well explain their own block.
- (ii) Block components are as correlated as possible for connected blocks.

$$\operatorname{cov}^{2}(\mathbf{X}_{j}\mathbf{w}_{j}, \mathbf{X}_{k}\mathbf{w}_{k}) = \operatorname{var}(\mathbf{X}_{j}\mathbf{w}_{j})\operatorname{cor}^{2}(\mathbf{X}_{j}\mathbf{w}_{j}, \mathbf{X}_{k}\mathbf{w}_{k})\operatorname{var}(\mathbf{X}_{k}\mathbf{w}_{k})$$

RGCCA for multigroup data analysis



Group loadings and group components should verify the following properties at the same time:

- Group component $\mathbf{X}_i \mathbf{w}_i$ well explains their own block.
- Small angle between group loadings if groups are connected.

 $\langle \mathbf{X}_{i}^{\mathsf{T}}\mathbf{X}_{i}\mathbf{w}_{i}, \mathbf{X}_{l}^{\mathsf{T}}\mathbf{X}_{l}\mathbf{w}_{l} \rangle = \cos(\mathbf{X}_{i}^{\mathsf{T}}\mathbf{X}_{i}\mathbf{w}_{i}, \mathbf{X}_{l}^{\mathsf{T}}\mathbf{X}_{l}\mathbf{w}_{l}) \times \|\mathbf{X}_{i}^{\mathsf{T}}\mathbf{X}_{i}\mathbf{w}_{i}\| \times \|\mathbf{X}_{l}^{\mathsf{T}}\mathbf{X}_{l}\mathbf{w}_{l}\|$

Multi-group data analysis The public good dilemma game

The Public Goods Game



Example from Kim De Roover et al.(2013) for illustrating Clusterwise SCA-ECP

Hypothesis

- 1. When both players contribute equally, they want to be rewarded equally.
- 2. Violation of equality is perceived as unjust and elicits anger.
- 3. In case of inequality, the subject's negative emotional reactions may be modified if the advantaged person is in personal need, because feelings of sympathy and empathy can be aroused.

An experimental study on 282 cooperators: Two factors "Reward" and "Empathy" are fully crossed



Reward is equal for both players (Equal condition). **Reward is higher** for the free-rider (Unequal condition).

High empathy condition: the cooperator receives a message about a negative personal event that happened to the free-rider. **Low empathy condition**: the same message is received, but the cooperator is asked to take an objective perspective. **Control condition**: the cooperator does not receive any message.

The public good dilemma game data set



Low Empathy High Empathy

Control



Clustering of groups according to cosine between grouploadings



PCA + Varimax on standardized groups



Technical conclusions

- 1. Depending on the dimension of the blocks, use either the **primal or the dual algorithm**.
- 2. The dual representation of the RGCCA algorithm allows:
 - Analysing high dimensional blocks.
 - Recovering **nonlinear relationships between blocks** (choice of the kernel function).
 - Handling any type of data (e.g. histogram) as long as relevant kernel is defined.

3. Sparse constraints are useful when the relevant variables are masked by (too many) noisy variables.

4. Kronecker constraints are used for threeway data configuration.

General conclusion



RGCCA for multiblock, multigroup or multiway data allows analyzing the data in their natural (but complex) structure.



Economic inequality and political instability Data from Russett (1964)

Economic inequality

Agricultural inequality

- **GINI :** Inequality of land distributions
- **FARM :** % farmers that own half of the land (> 50)
- **RENT :** % farmers that rent all their land

Industrial development

- **GNPR :** Gross national product per capita (\$ 1955)
- LABO : % of labor force employed in agriculture

Political instability

- **INST :** Instability of executive (45-61)
- ECKS : Nb of violent internal war incidents (46-61)
- **DEAT :** Nb of people killed as a result of civic group violence (50-62)
- **D-STAB :** Stable democracy
- **D-UNST :** Unstable democracy
- **DICT**: Dictatorship

Economic inequality and political instability (Data from Russett, 1964)



Agricultural inequality

GINI : Inequality of land distributions

FARM : % farmers that own half of the land (> 50)

RENT : % farmers that rent all their land

Industrial development

GNPR : Gross national product per capita (\$, 1955) **LABO :** % of labor force employed in agriculture

Political instability

INST : Instability of executive (45-61)
ECKS : Nb of violent internal war incidents (46-61)
DEAT : Nb of people killed as a result of civic group violence (50-62)
DEMO : Stable democracy (1), Unstable democracy (2) or Dictatorship (3)

Path Diagram

Agricultural inequality (X₁)



Outer weight vectors

Agricultural inequality (X₁)



Industrial development (X₂)

Political instability (X₃)

Bootstrap confidence intervals





Data vizualization

