# Group-wise Principal Component Analysis

José Camacho, Edoardo Saccenti, Roberto Therón

SSC¹⁵
NAANTALI • FINLAND
19-22 June 2017

15th Scandinavian Symposium on Chemometrics

Network Engineering & Security Group
http://nesg.ugr.es

ugr | Universidad de Granada

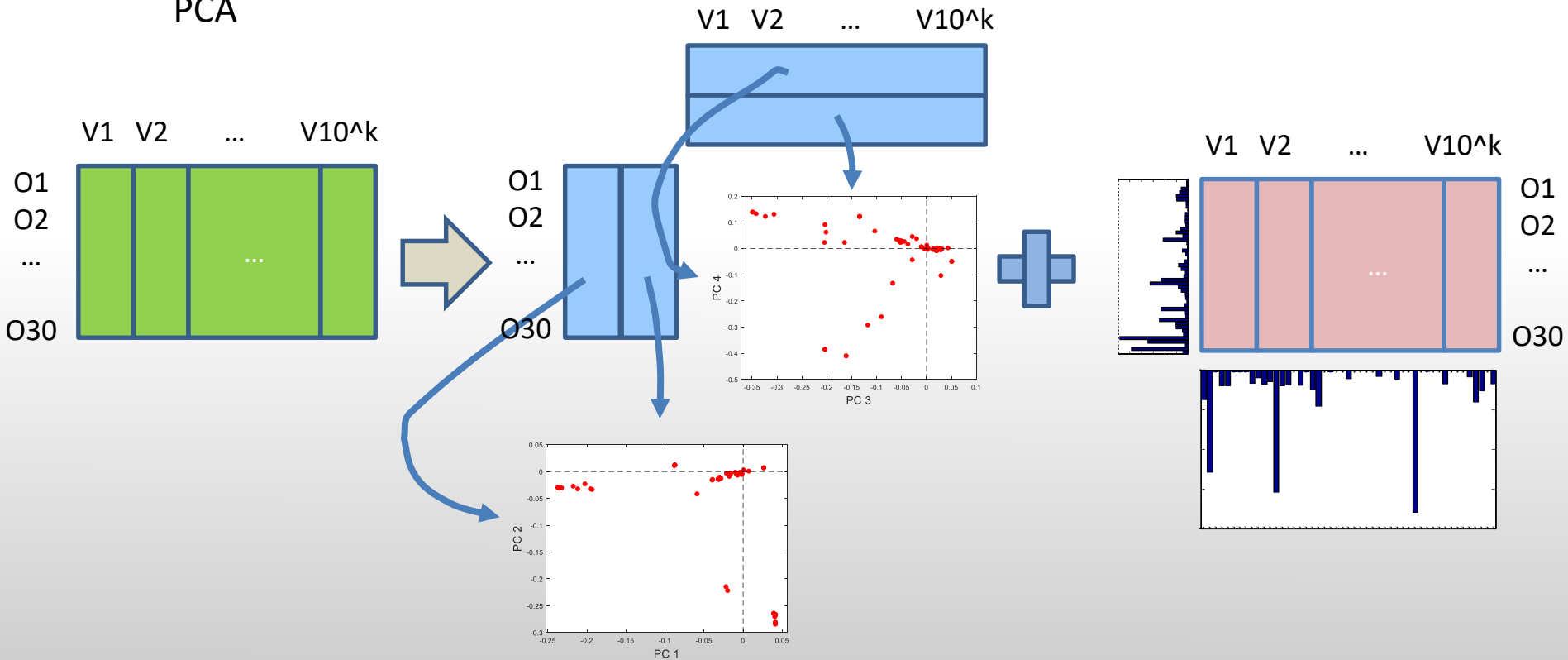# ➡ Exploratory Data Analysis (EDA)

✔ ***(Human) Learning from data****: to improve the understanding of a phenomenon of interest by analyzing data collected on a number of (hopefully) relevant variables.*
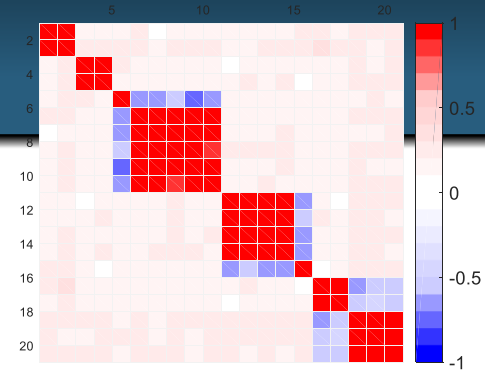
- *Statistics*
- *Visual Analytics*
- *Machine learning*
- *...*

➡ Multivariate EDA approach: Matrix Factorization

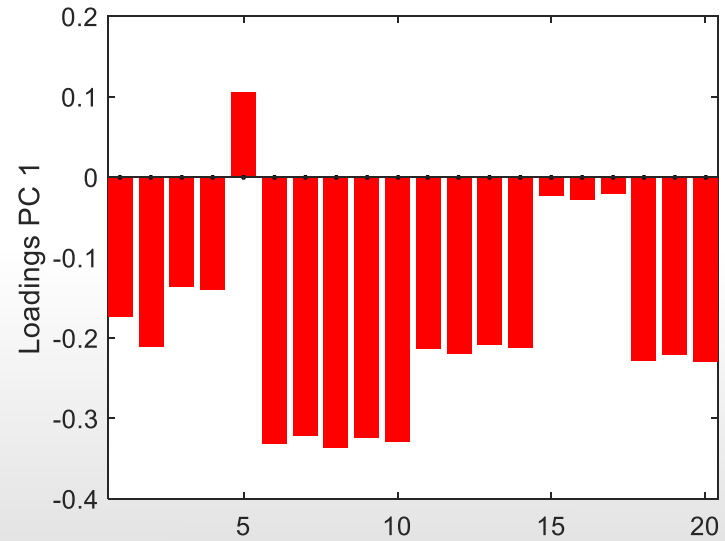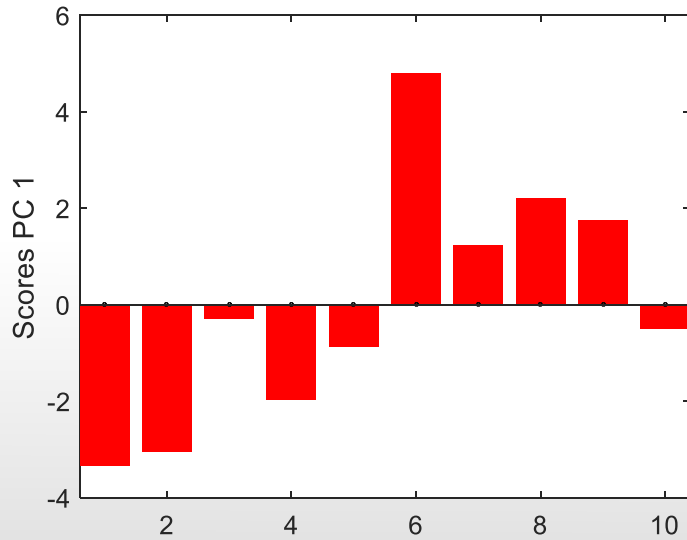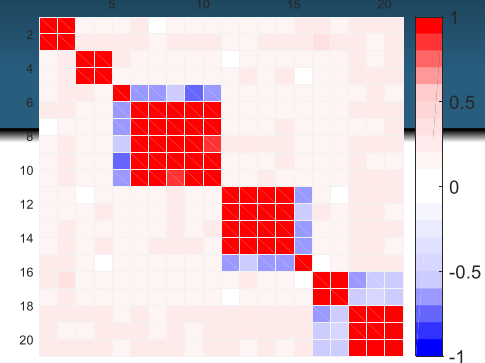$$\mathbf{X} = \mathbf{T}_A \cdot \mathbf{P}_A^T + \mathbf{E}_A$$

→ PCA

✔ Structure ≈ Maximum variance

✔ PCA for EDA? ➔ X(20x10) = [[1:2], [3:4], [5:10], [11:15], [16:20]]
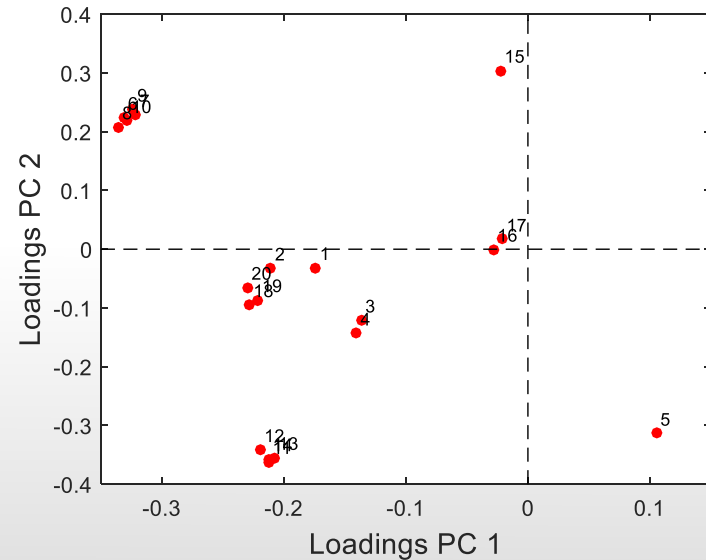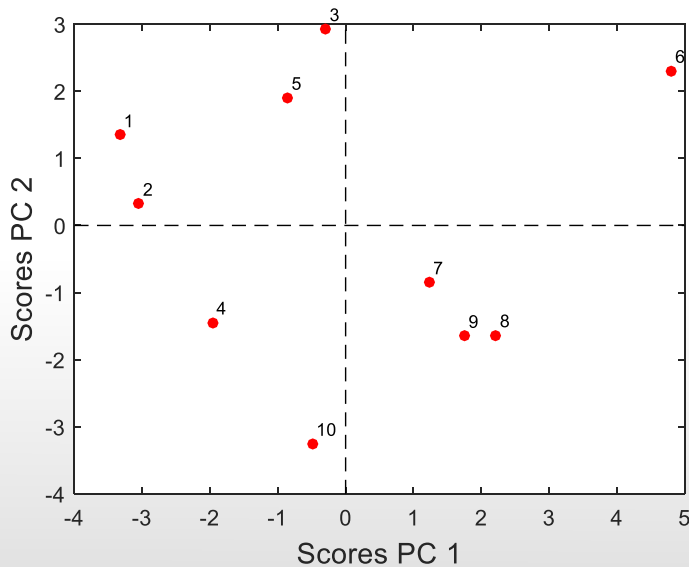


First PC

➡ PCA

   ✔ Structure ≈ Maximum variance
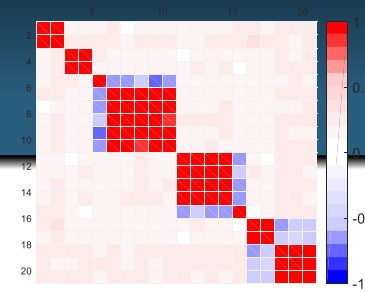
   ✔ PCA for EDA? ➡ X(20x10) = [[1:2], [3:4], [5:10], [11:15], [16:20]]
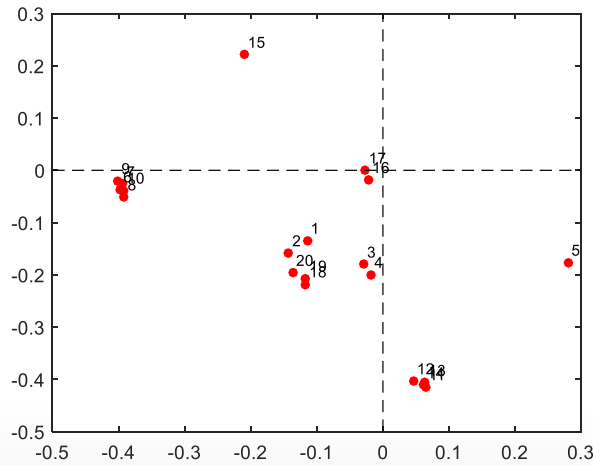
First 2 PCs

➡ PCA ⬅ ? ➡ Relationship among Variables

 ☑ Jackson, Jolliffe ➡ NO

 ☑ PCA does not distinguish between unique variance and shared variance

   ✓ Factor Analysis ➡ Model Shared Variance

 ☑ The PCA factorization is poorly interpretable because the principal
   components are linear combinations of all the variables

   ✓ Rotation

   ✓ Sparse Methods        ➡ Trade-off between variance and simplicity

 ☑ 1 PC contains many SoV & 1 SoV in many PCs ➡ GPCA (without biplots)
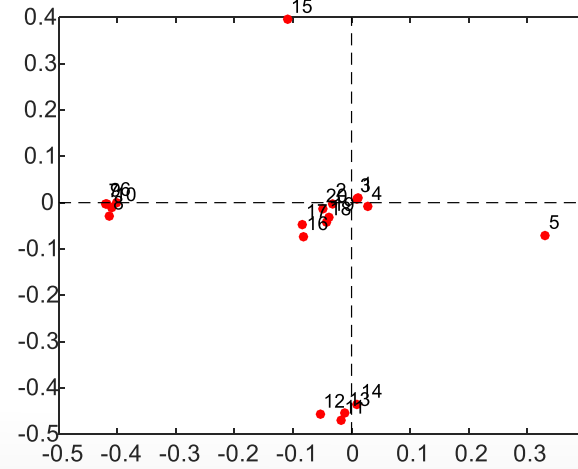
http://nesg.ugr.es

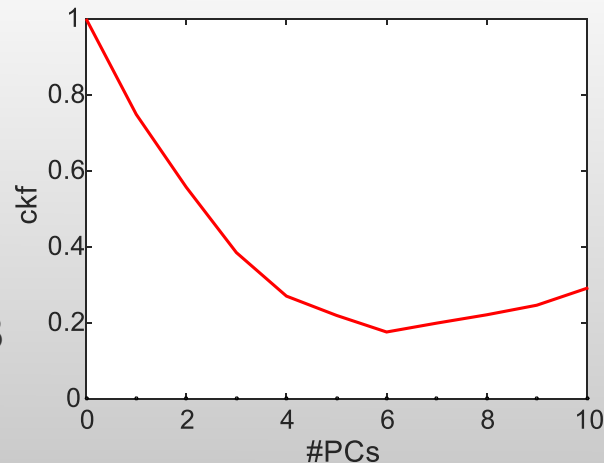➡ Variance vs Simplicity:

✔ PCA+ Varimax ➡ Rotation depends on #PCs (and scaling)
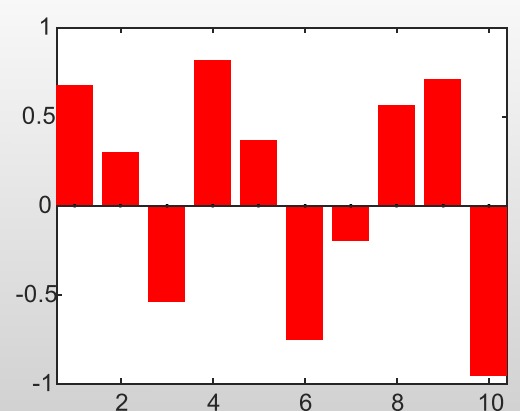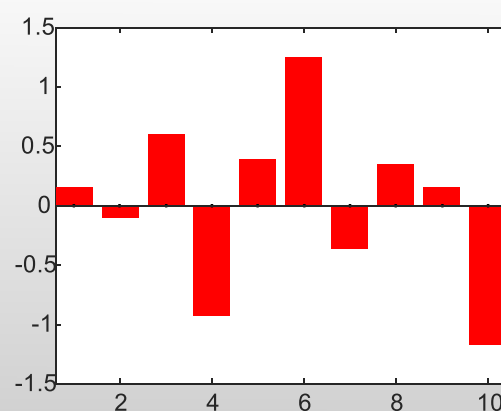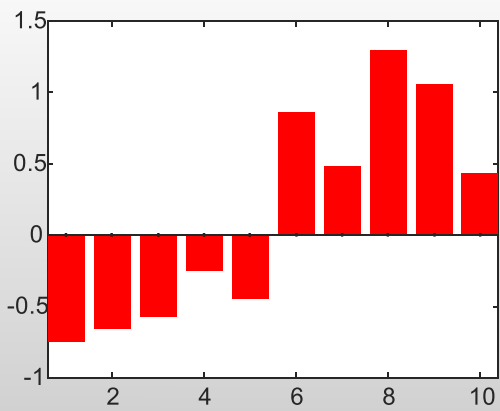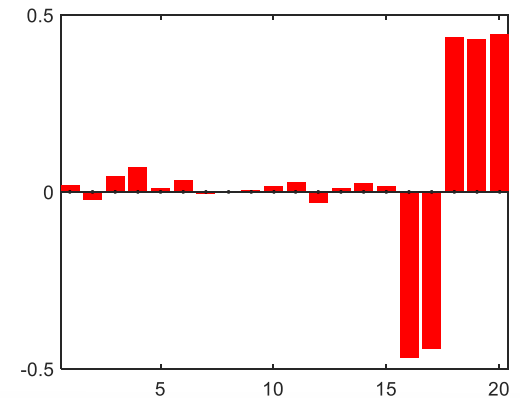


Rotate from 2 PCs

Rotate from 6 PCs
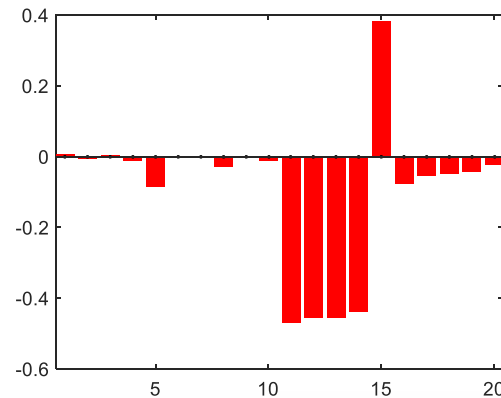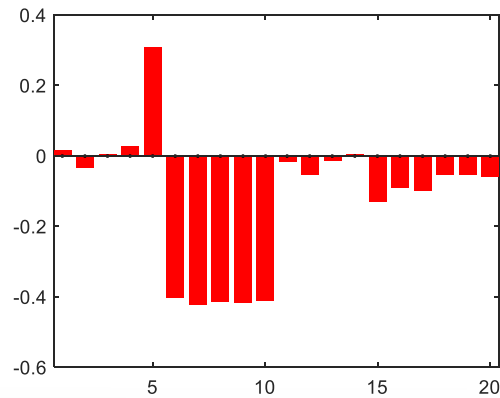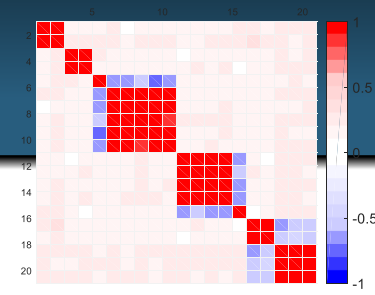
Fast PCA CV (ckf)
J. Chem. 29(2015):467–478

Variance vs Simplicity:

✔ PCA+ Varimax: Rotate from 6 PCs

http://nesg.ugr.es

➡ Variance vs Simplicity:

✔ SPCA: Depends on metaparameters



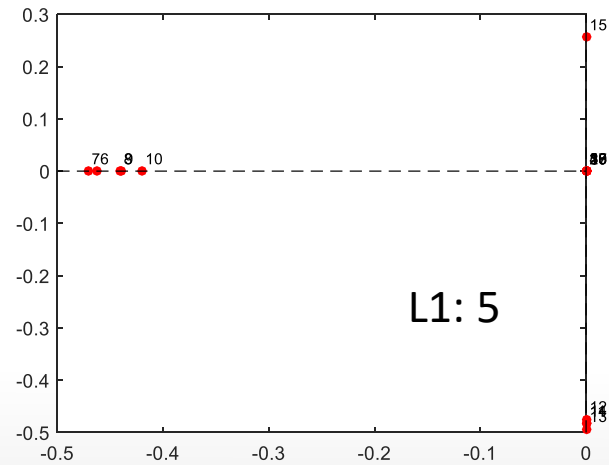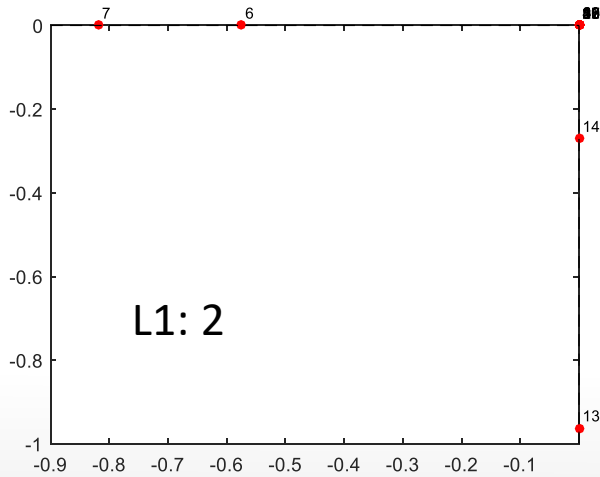ckf



**Group-wise Principal Component Analysis (GPCA)**

http://nesg.ugr.es

➡ Variance vs Simplicity:

✔ PCA+ Varimax: Rotate from 6 PCs

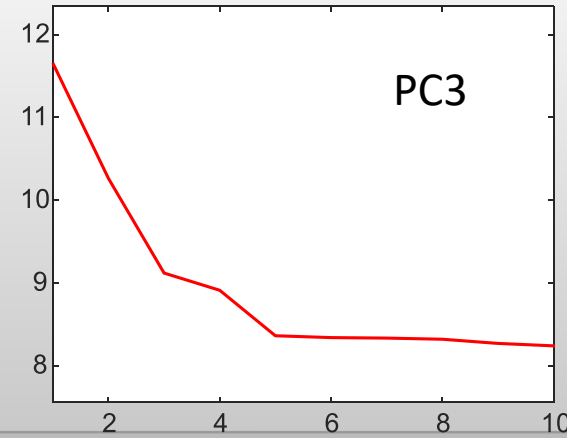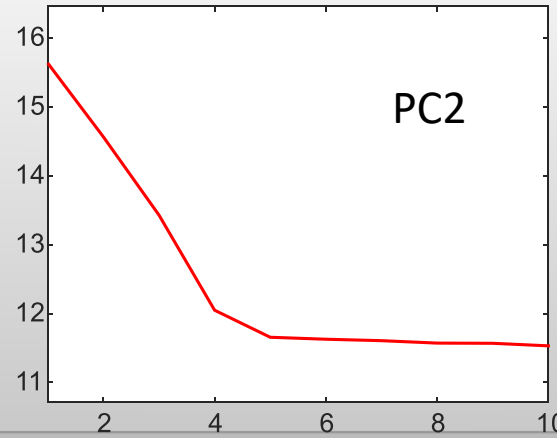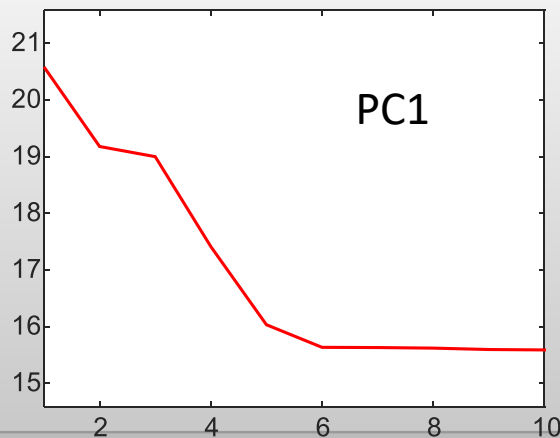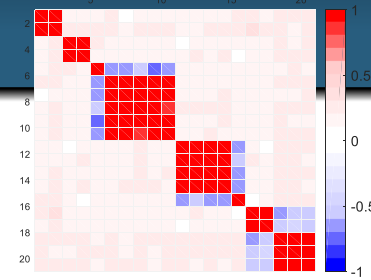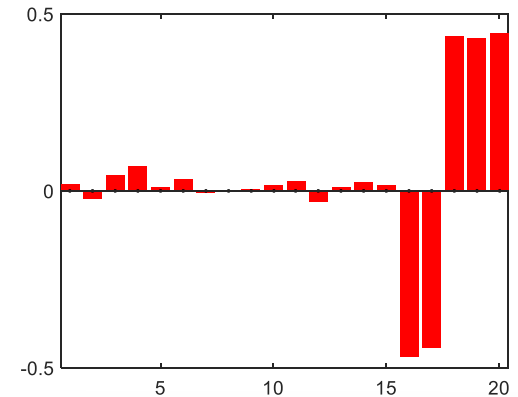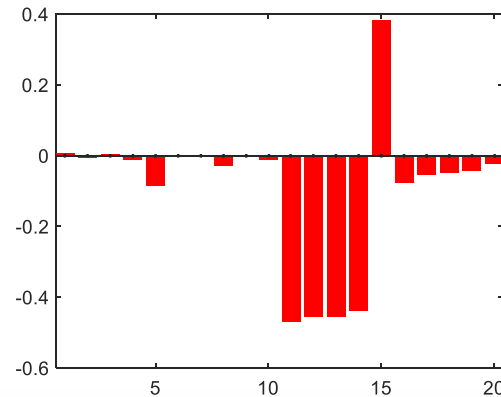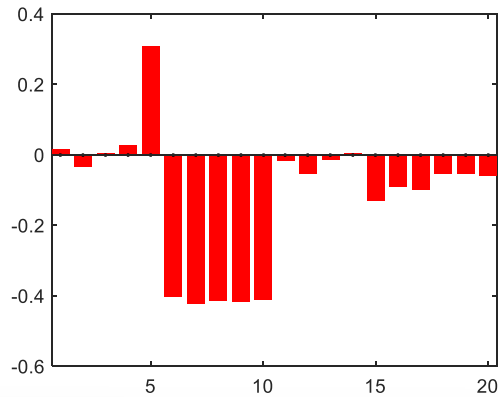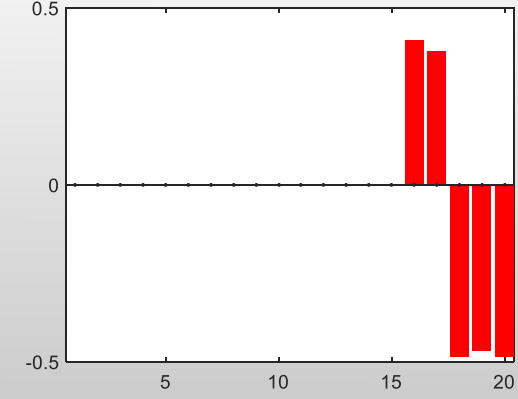

✔ SPCA: L1: [6,5,5]

➡️ Variance vs Simplicity:

- ✔ Result depends on a good choice of metaparameters
- ✔ Typical approach (regression) inherited ➜ CV
- ✔ Risks
  - Prediction ≠ Interpretation
    - Oversimplify / Overcomplicate
    - Application to non-sparse data
  - PCA CV ≠ PLS CV (problem with independent variables, see Journal of Chemometrics, 2012, 26 (7): 361-373.)

# ➡ Missing-Data for EDA (MEDA)

## Missing-data theory in the context of exploratory data analysis
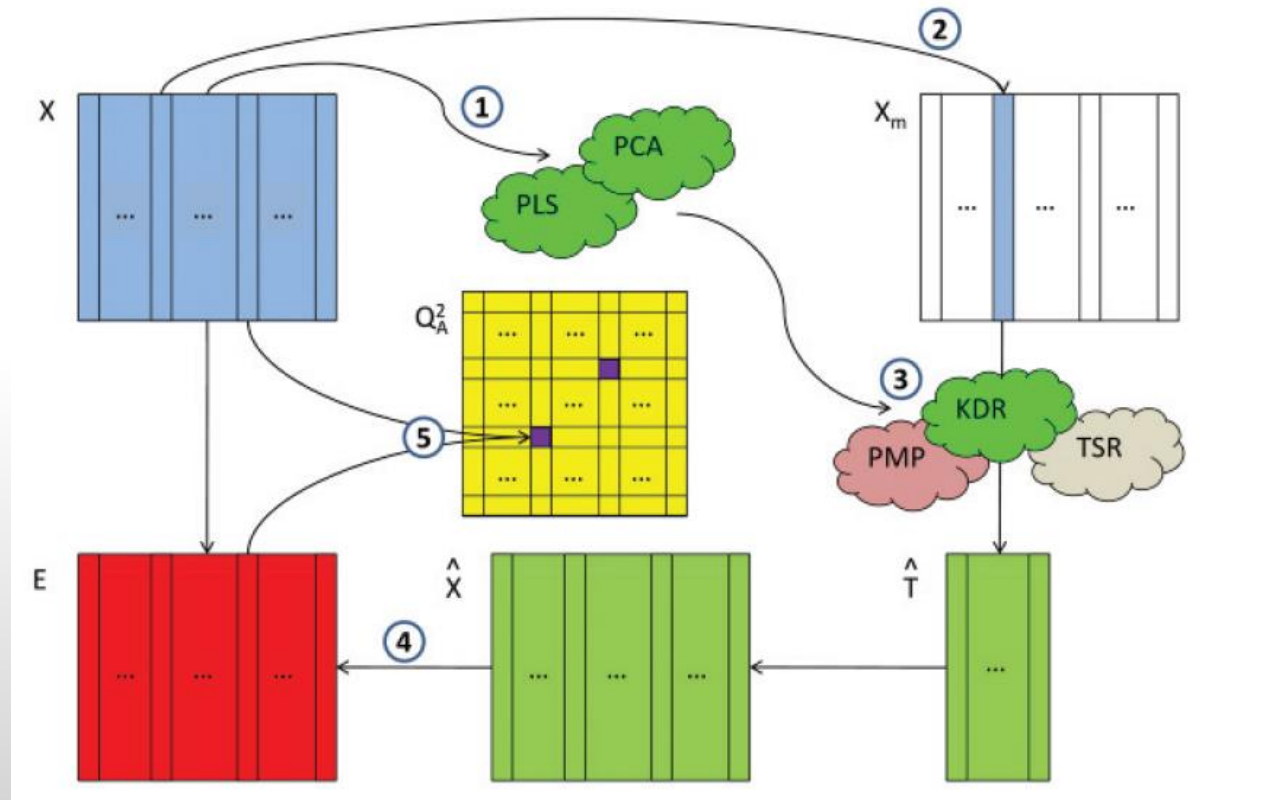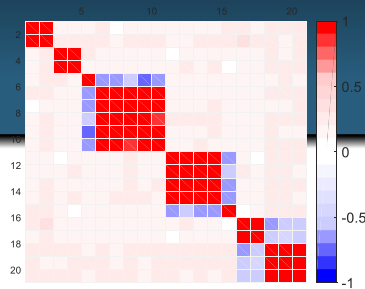
José Camacho

*Departamento de Teoría de la Señal, Telemática y Comunicaciones, Universidad de Granada, 18071, Granada, Spain*
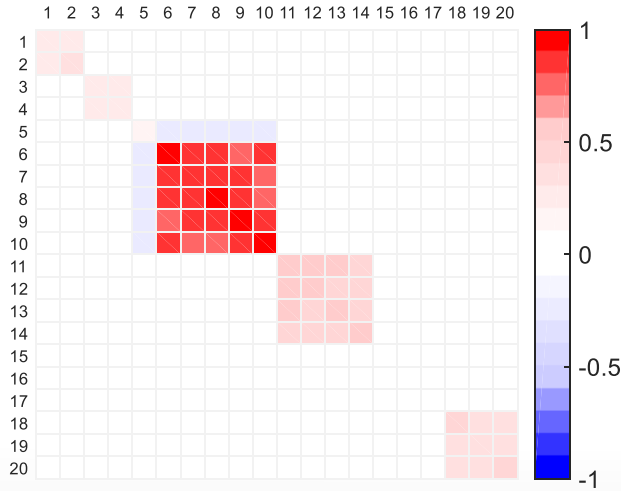
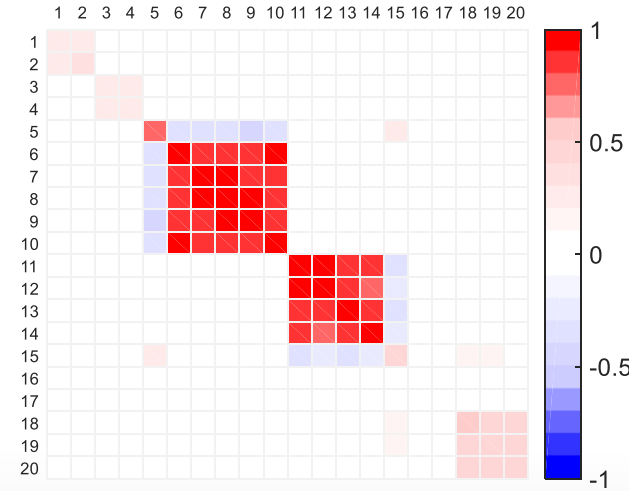Instead of changing the model, change the visualization
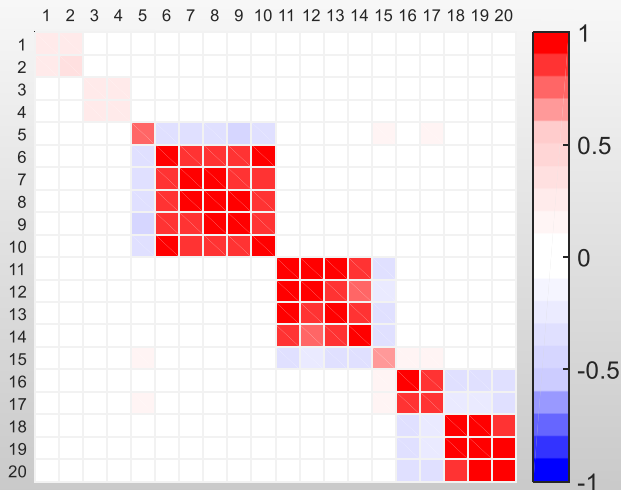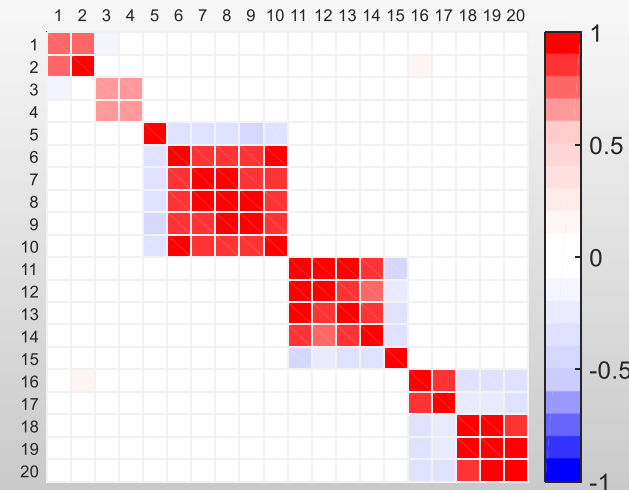
# Missing-Data for EDA (MEDA)

## MEDA:



1 PC

1-2 PCs

1-3 PCs

1-4 PCs

**Group-wise Principal Component Analysis (GPCA)**

# ➡ Group-wise PCA:

## Group-Wise Principal Component Analysis for Exploratory Data Analysis

José Camacho [a], Rafael A. Rodríguez-Gómez [a], and Edoardo Saccenti [b]

[a]Department of Signal Theory, Networking and Communication, University of Granada, Granada, Spain; [b]Laboratory of Systems and Synthetic Biology, Wageningen University & Research Center, Wageningen, The Netherlands

X   Do not force a data-driven simple structure

✓   Find the structure and force it in the model

1 PC ⬅➡ 1 SoV

# Group-wise PCA:

## Three steps:

1. Find structure (MEDA)

2. Identify Groups of Variables (Group Identification Algorithm or GIA)

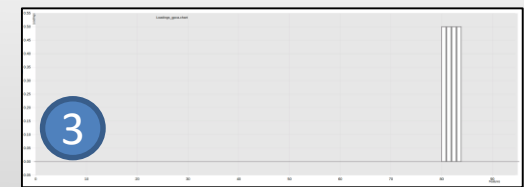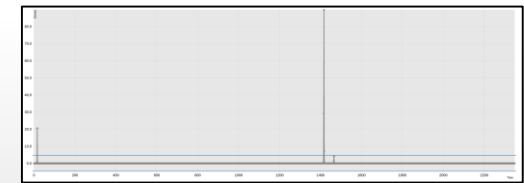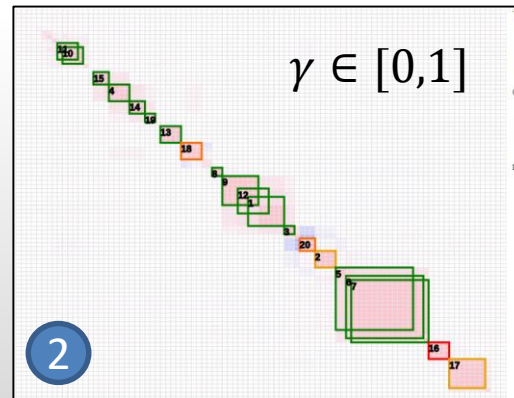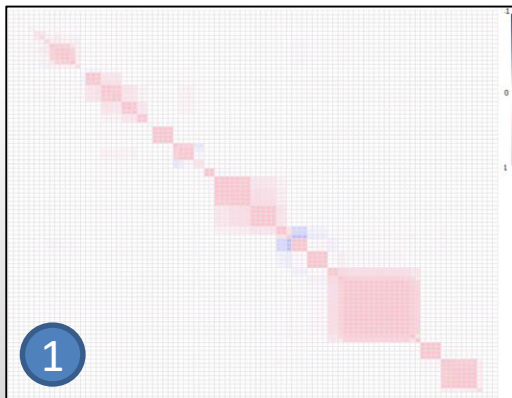3. Calibrate a group-wise PCA model (GPCA)



$\gamma \in [0,1]$

# ➡ Group-wise PCA:

✔ Initialize: $\quad C = X^T X$

$$B = I,$$

✔ For each PC

- For each ($k$-th) group in GIA $\quad C^k = C$

$$c^k_{lm} = 0, \ \forall l \notin S_k \ or \ \forall m \notin S_k.$$

- Compute 1 PC: $\quad C^k = p^k (\sigma^k)^2 (p^k)^T + E^k.$

✔ Choose PC with most variance:

$$p_a = \arg\min_{p^k} \|E^k\|_F$$

$$t_a = X p_a.$$

✔ Deflate (Mackey'08):

$$q = B p_a$$
$$C = (I - qq^T) C (I - qq^T)$$
$$B = B(I - qq^T).$$

GPCA: X(20x10) = [1, 2, 3, [4:5], [6:9], 10]

$\gamma = 0.2$

Visually selected

http://nesg.ugr.es

➡ GPCA: X(20x10) = [1, 2, 3, [4:5], [6:9], 10]

$\gamma = 0.2$

Visually selected



Group-wise Principal Component Analysis (GPCA)

# Want to play with GPCA?

## iGPCA Dashboard 1.0

*iGPCA, the interactive GPCA analysis.*

See J. Camacho, R. A. Rodríguez-Gómez, and E. Saccenti, "Group-wise Principal Component Analysis for Exploratory Data Analysis," Journal of Computational and Graphical Statistics, pp. 0–0, Dec. 2016. for more details.

START ANALYSIS    RESET

http://nesg.ugr.es:5003

**MEDA**

**Multivariate Exploratory Data Analysis**

| VERSION 1.0 January 2015 | José Camacho Páez Rafael Rodríguez Gómez Alejandro Pérez Villegas Elena Jiménez Mañas |

PCA    PLS

ChemoLab, (2015) 143: 49

https://github.com/josecamachop/MEDA-Toolbox

➡ Multivariate Exploratory Data Analysis can be tricky

- ✔ Variance vs Simplicity selected by prediction (CV)
- ✔ Rather: find structure in a EDA manner and impose it in model.

➡ GPCA:

- ✔ It is sparse when data is group-wise (in the variables)
- ✔ Only correlated variables (1 SoV) in a PC
- ✔ Does not oversimplify/overcomplicate structure
- ✔ Metaparameter selected from visualization (perfect for EDA)
- ✔ Still, you can always set GPCA by CV

# Group-wise Principal Component Analysis

**José Camacho, Edoardo Saccenti, Roberto Therón**

This work is partly supported by the Spanish Ministry of Economy and Competitiveness and FEDER funds through project TIN2014-60346-R

GOBIERNO DE ESPAÑA
MINISTERIO DE ECONOMÍA Y COMPETITIVIDAD

**Unión Europea**
Fondo Europeo de Desarrollo Regional

Network Engineering & Security Group
http://nesg.ugr.es

ugr | Universidad de Granada