

# The RGCCA package for Regularized/Sparse Generalized Canonical Correlation Analysis

Etienne CAMENEN

2019-12-18

## Contents

<b>1</b>	<b>Multiblock data analysis with the RGCCA package</b>	<b>1</b>
<b>2</b>	<b>Load the inputs</b>	<b>1</b>
2.1	Load the blocks . . . . .	1
2.2	Load the groups of response and the connection between blocks . . . . .	1
2.3	View the inputs . . . . .	2
<b>3</b>	<b>Run S/RGCCA</b>	<b>3</b>
<b>4</b>	<b>Vizualise the analysis</b>	<b>3</b>
4.1	With the superbloc, by default on the first and the second components . . . . .	4
4.2	With the politic block, on the 2nd and the 3rd components . . . . .	6

## 1 Multiblock data analysis with the RGCCA package

We consider  $J$  data matrices  $X_1, \dots, X_J$ . Each  $n \times p_j$  data matrix  $X_j = [x_{j1}, \dots, x_{jp_j}]$  is called a block and represents a set of  $p_j$  variables observed on  $n$  individuals. The number and the nature of the variables may differ from one block to another, but the individuals must be the same across blocks. We assume that all variables are centered. The objective of RGCCA is to find, for each block, a weighted composite of variables (called block component)  $y_j = X_j \cdot a_j$ ,  $j = 1, \dots, J$  (where  $a_j$  is a column-vector with  $p_j$  elements) summarizing the relevant information between and within the blocks. The block components are obtained such that (i) block components explain well their own block and/or (ii) block components that are assumed to be connected are highly correlated. In addition, RGCCA integrates a variable selection procedure, called SGCCA, allowing the identification of the most relevant features.

## 2 Load the inputs

### 2.1 Load the blocks

The blocks are loaded with the function `load_blocks`. The first argument of this function (`superblock`) required a boolean giving the presence (TRUE) / absence (FALSE) of a superbloc. The second one corresponds to a character giving the list of the file path separated by a comma (argument `file`). By default, the name of the blocks corresponds to those of the files (`names` argument) and could be set. By default, the tabulation is used as a column separator (`sep` argument) and the first row is considered as a header (`header` parameter).

```
# Warning : separators by default are tabulation
blocks = load_blocks(file = "data/agriculture.tsv,data/industry.tsv,data/politic.tsv")
```

### 2.2 Load the groups of response and the connection between blocks

The connection between the blocks will be used by the RGCCA and must be set by `set_connection` function. A group of samples will be used to color them in the samples plot and must be set by `load_response` function. For both functions, the `blocks` parameter, set at the previous step, is required. The other parameters are optional. The user

could import a file containing either (`file` parameter) : (i) a symmetric matrix with 1 giving a connection between two blocs, or 0 otherwise; (ii) a univariate vector (qualitative or quantitative) or a disjunctive table for the response. By default, the column separator is the tabulation and could be set (`sep` argument). For the `load_response`, a header could be specified (`header` parameter).

```
# Optional parameters
response <- connection <- NULL

# Uncomment the parameters below to try without this settings
response = load_response(blocks = blocks, file = "data/response.tsv")
connection = load_connection(file = "data/connection.tsv")
```

## 2.3 View the inputs

Table 1: agriculture

	gini	farm	rent
<b>Argentina</b>	86.3	98.2	3.52
<b>Australia</b>	92.9	99.6	3.27
<b>Austria</b>	74	97.4	2.46
<b>Belgium</b>	58.7	85.8	4.15
<b>Bolivia</b>	93.8	97.7	3.04
<b>Brasil</b>	83.7	98.5	2.31

Table 2: industry

	gnpr	labo
<b>Argentina</b>	5.92	3.22
<b>Australia</b>	7.1	2.64
<b>Austria</b>	6.28	3.47
<b>Belgium</b>	6.92	2.3
<b>Bolivia</b>	4.19	4.28
<b>Brasil</b>	5.57	4.11

Table 3: politic

	inst	ecks	death	demostab	demoinst	dictator
<b>Argentina</b>	0.07	4.06	5.38	0	1	0
<b>Australia</b>	0.01	0	0	1	0	0
<b>Austria</b>	0.03	1.61	0	0	1	0
<b>Belgium</b>	0.45	2.2	0.69	1	0	0
<b>Bolivia</b>	0.37	3.99	6.5	0	0	1
<b>Brasil</b>	0.45	3.91	0.69	0	1	0

Table 4: response

<b>Argentina</b>	demoinst
<b>Australia</b>	demostab
<b>Austria</b>	demoinst

<b>Belgium</b>	demostab
<b>Bolivia</b>	dictator
<b>Brasil</b>	demoinst

Table 5: connection

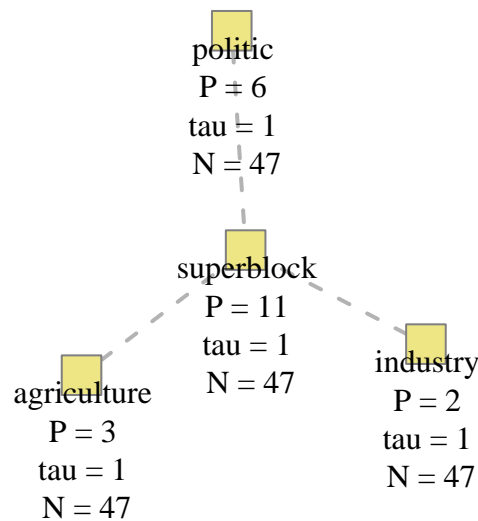
0	0	0	1
0	0	0	1
0	0	0	1
1	1	1	0

### 3 Run S/RGCCA

SGCCA is run from the RGCCA package by using two components for a biplot visualization. The S/RGCCA function doesn't names the blocks in their outputs. This step is required to generate biplots.

```
sgcca_out = rgcca.analyze(blocks = blocks)
plot_network(sgcca_out)
```

**Common rows between blocks : 47**



### 4 Vizualise the analysis

Both the samples and the variables could be visualized by using biplots functions (respectively `plot_ind` and `plot_var_2D`). Histograms are used to visualized in decreasing order the variables with the higher weights and the blocks with the higher Average Variance Explained (AVE).

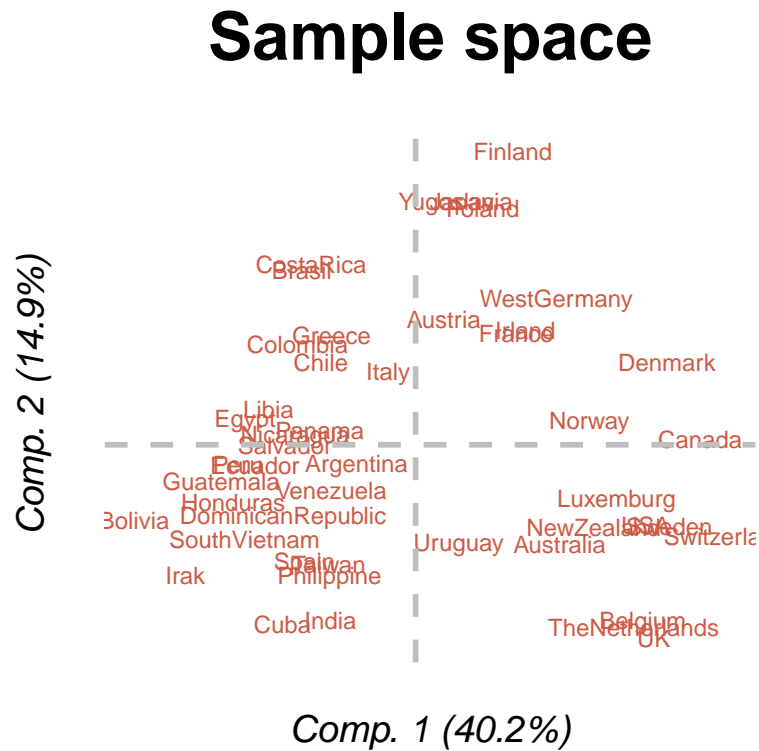
These functions take the results of a sgcca or a rgcca (`rgcca` parameter) and the components to visualize : either `compx` and `compy` for biplots or `comp` for histograms. By default, `compx = comp = 1` and `compy = 2`. The presence or the absence of a superblock among the analysis could be specified for `plot_var_2D` and `plot_var_1D` to color the variables according to their blocks. By default, the last block is plotted, corresponding to the superblock if selected (`i_block` parameter). `plot_var_2D`, which is a corcircle plot, required the `blocks` for the correlation with the selected component. `plot_var_2D` could use the response variable to color the samples by groups. By default, the first 100th higher weights are used for the `plot_var_1D` and could be set by using the `n_mark` argument.

```
comp1 = 1
comp2 = 2
nmark = 100
```

## 4.1 With the superblock, by default on the first and the second components

### 4.1.1 Samples plot

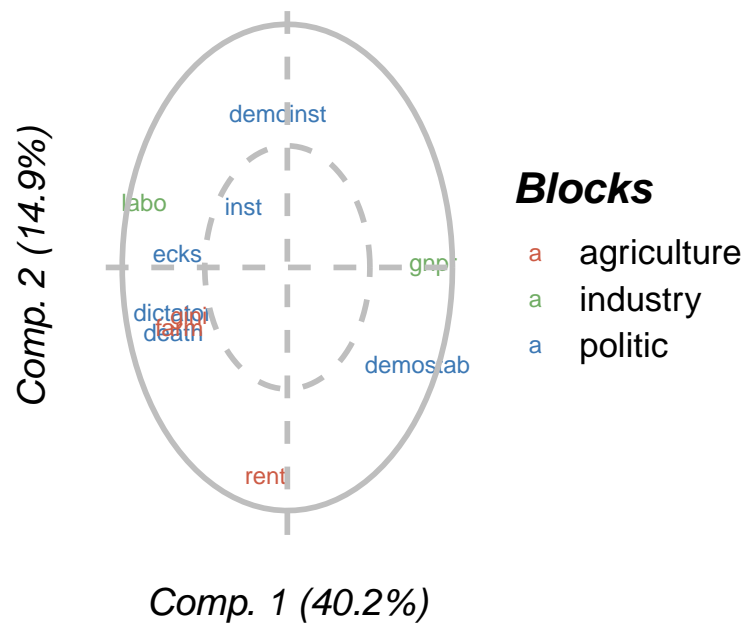
```
plot_ind(sgccca_out)
```



### 4.1.2 Corcircle plot

```
plot_var_2D(sgccca_out)
```

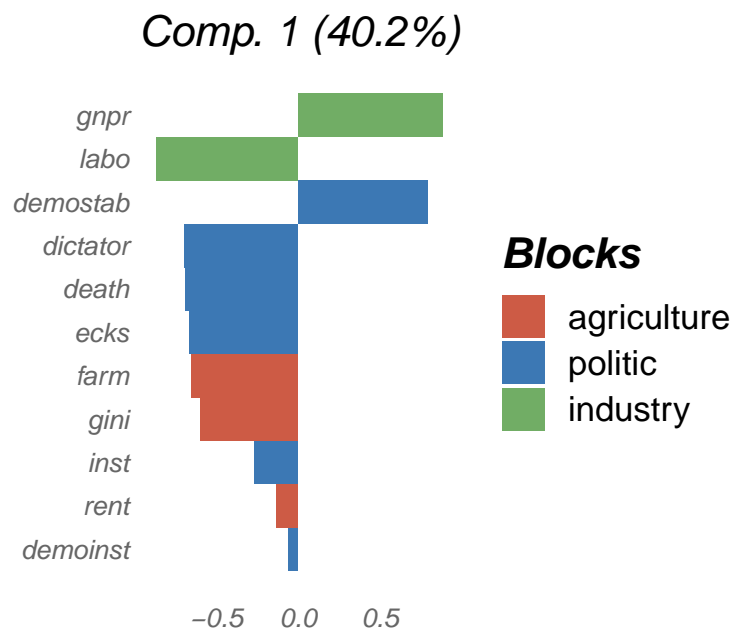
# Variable space



## 4.1.3 Fingerprint plot

`plot_var_1D(sgccca_out)`

# riable correlations with

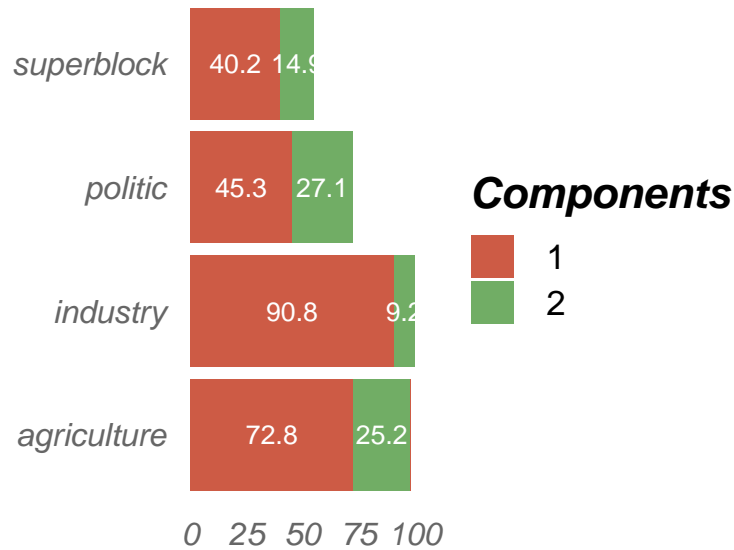


#### 4.1.4 Best explained blocks

```
plot_ave(sgccca_out)
```

## Age Variance Explained

*First outer comp. : 50.6% & 19.1%*

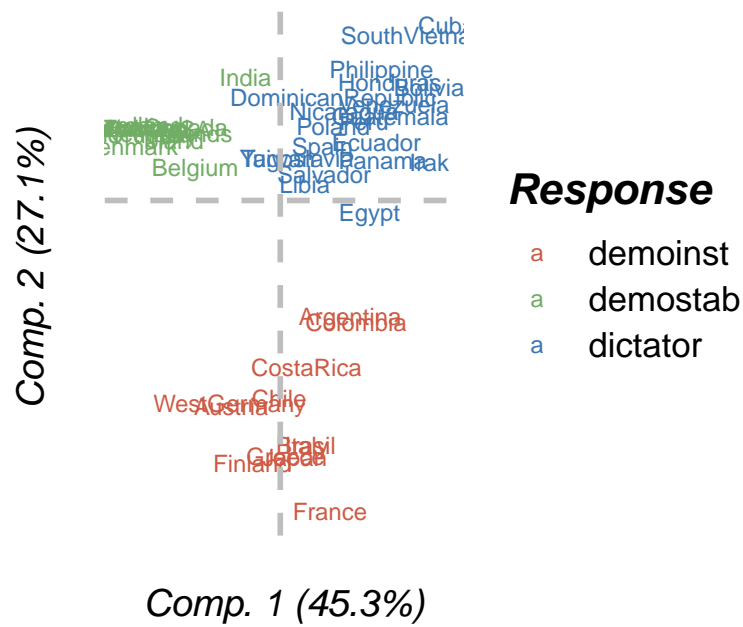


## 4.2 With the politic block, on the 2nd and the 3rd components

### 4.2.1 Samples plot

```
plot_ind(  
  rgcca = sgcca_out,  
  resp = response,  
  compx = comp1,  
  compy = comp2,  
  i_block = 3)
```

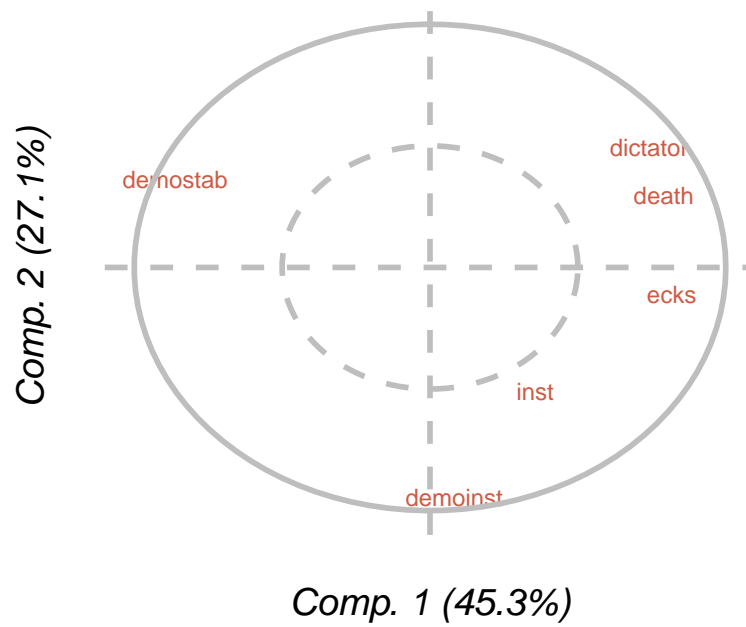
# Sample space



## 4.2.2 Corcircle plot

```
plot_var_2D(
  rgcca = sgcca_out,
  compx = comp1,
  compy = comp2,
  i_block = 3)
```

# Variable space



## 4.2.3 Fingerprint plot

```
plot_var_1D(  
    rgcca = sgcca_out,  
    comp = comp1,  
    n_mark = nmark,  
    i_block = 3,  
    type = "weight")
```



# Variable weights on

Comp. 1 (45.3%)

