



SOMbrero: an R package for numeric and non-numeric self-organizing maps

Nathalie Villa-Vialaneix
with J. Boelaert, L. Bendhaïba, M. Olteanu

nathalie.villa@toulouse.inra.fr

<http://www.nathalievilla.org>

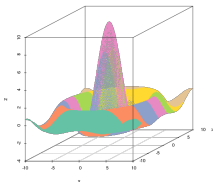


WSOM 2014 - Mittweida, Germany - July 4th



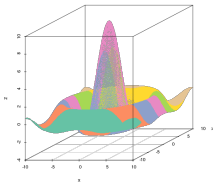


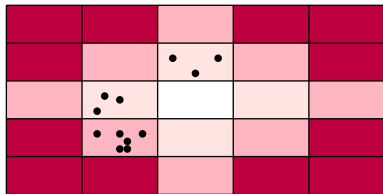
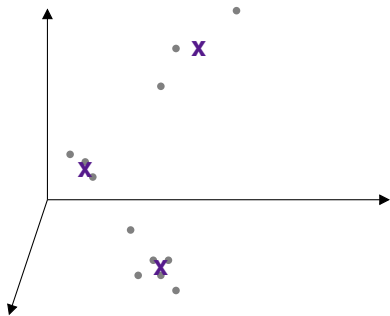
- 1 a short review of Self-Organizing Maps for non vectorial data
- 2 **SOMbrero**





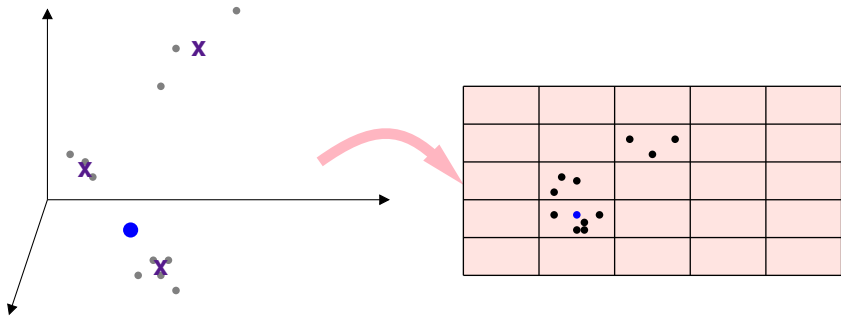
- 1 a short review of Self-Organizing Maps for non vectorial data
- 2 SOMbrero





- $(x_i)_{i=1,\dots,n} \subset \mathbb{R}^d$ are affected to a unit $C(x_i) \in \{1, \dots, U\}$
- the grid is equipped with a “distance” between units: $d(u, u')$ and observations affected to close units are close in \mathbb{R}^d
- every unit u corresponds to a **prototype**, $p_u(\mathbf{x})$ in \mathbb{R}^d

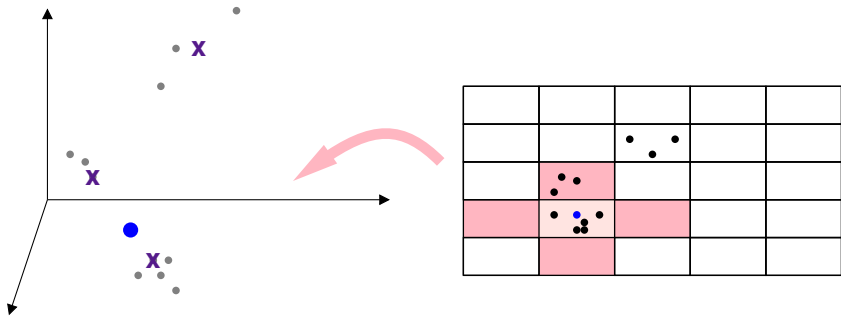




Iterative learning (affectation step): x_i is picked at random within $(x_k)_k$ and affected to *best matching unit*:

$$C(x_i) = \arg \min_u \|x_i - p_u\|^2$$





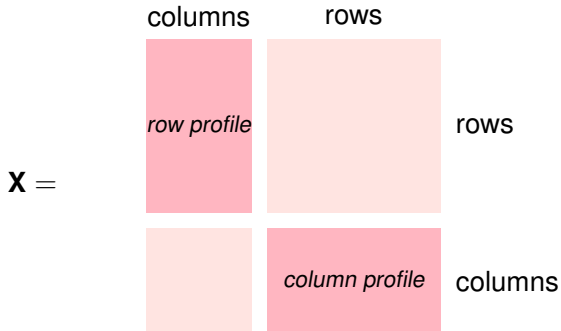
Iterative learning (representation step): all prototypes in neighboring units are updated with a gradient descent like step:

$$p_u^{t+1} \leftarrow p_u^t + \mu(t) H^t(d(C(x_i), u))(x_i - p_u^t)$$





Data: contingency table $\mathbf{T} = (n_{ij})_{ij}$ with p rows and q columns transformed into a numeric dataset \mathbf{X} :



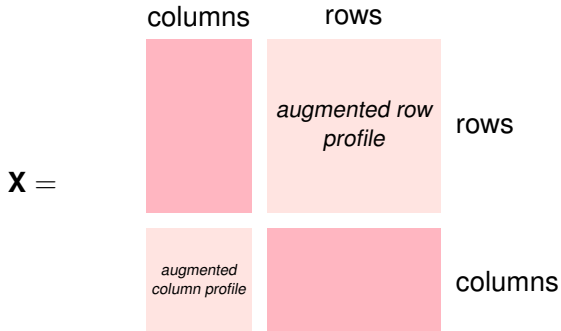
with

- $\forall i = 1, \dots, p$ and $\forall j = 1, \dots, q$, $\mathbf{x}_{ij} = \frac{n_{ij}}{n_i} \times \sqrt{\frac{n}{n_j}}$





Data: contingency table $\mathbf{T} = (n_{ij})_{ij}$ with p rows and q columns transformed into a numeric dataset \mathbf{X} :



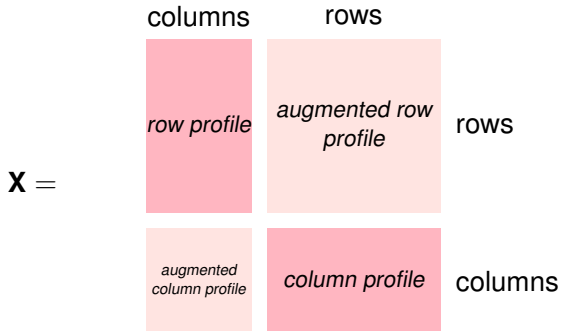
with

- $\forall i = 1, \dots, p$ and $\forall j = q + 1, \dots, q + p$, $\mathbf{x}_{ij} = \mathbf{x}_{k(i)+p,j}$ with $k(i) = \arg \max_{k=1, \dots, q} \mathbf{x}_{ik}$





Data: contingency table $\mathbf{T} = (n_{ij})_{ij}$ with p rows and q columns transformed into a numeric dataset \mathbf{X} :



- **affectation** uses reduced profile
- **representation** uses augmented profile
- alternatively process row profiles and column profiles





[[Hammer and Hasenfuss, 2010](#), [Olteanu and Villa-Vialaneix, 2014](#)]

Data: described by a dissimilarity matrix $\mathbf{D} = (\delta(x_i, x_j))_{i,j=1,\dots,n}$
($(x_j)_i$ not necessarily vectorial)





[[Hammer and Hasenfuss, 2010](#), [Olteanu and Villa-Vialaneix, 2014](#)]

Data: described by a dissimilarity matrix $\mathbf{D} = (\delta(x_i, x_j))_{i,j=1,\dots,n}$

$((x_i)_i$ not necessarily vectorial)

Adaptations of the SOM algorithm:

- **prototypes:** expressed as (symbolic) convex combination of $(x_i)_i$: $p_u \sim \sum_{i=1}^n \gamma_{ui} x_i$, $\gamma_{ui} \geq 0$ and $\sum_i \gamma_{ui} = 1$





[[Hammer and Hasenfuss, 2010](#), [Olteanu and Villa-Vialaneix, 2014](#)]

Data: described by a dissimilarity matrix $\mathbf{D} = (\delta(x_i, x_j))_{i,j=1,\dots,n}$

$((x_j)_i)$ not necessarily vectorial)

Adaptations of the SOM algorithm:

- **prototypes:** expressed as (symbolic) convex combination of $(x_i)_i$: $p_u \sim \sum_{i=1}^n \gamma_{ui} x_i$, $\gamma_{ui} \geq 0$ and $\sum_i \gamma_{ui} = 1$
- **distance computation:** $\|x_j - p_u\|^2$ replaced by

$$(\mathbf{D}\gamma_u)_i - \frac{1}{2}\gamma_u^T \mathbf{D}\gamma_u$$

in reference to a pseudo-Euclidean framework [[Goldfarb, 1984](#)]





[[Hammer and Hasenfuss, 2010](#), [Olteanu and Villa-Vialaneix, 2014](#)]

Data: described by a dissimilarity matrix $\mathbf{D} = (\delta(x_i, x_j))_{i,j=1,\dots,n}$

((x_j)_{*i*} not necessarily vectorial)

Adaptations of the SOM algorithm:

- **prototypes:** expressed as (symbolic) convex combination of (x_i)_{*i*}: $p_u \sim \sum_{i=1}^n \gamma_{ui} x_i$, $\gamma_{ui} \geq 0$ and $\sum_i \gamma_{ui} = 1$
- **distance computation:** $\|x_j - p_u\|^2$ replaced by

$$(\mathbf{D}\gamma_u)_i - \frac{1}{2}\gamma_u^T \mathbf{D}\gamma_u$$

in reference to a pseudo-Euclidean framework [[Goldfarb, 1984](#)]

- **representation:** replaced by an update of (γ_u)_{*u*}:

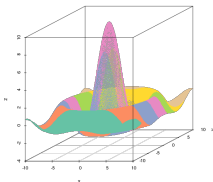
$$\gamma_u^{t+1} \leftarrow \gamma_u^t + \mu(t) H^t(d(C(x_i), u)) (\mathbf{1}_i - \gamma_u^t)$$

with $\mathbf{1}_{il} = 1$ if $l = i$ and 0 otherwise.





- 1 a short review of Self-Organizing Maps for non vectorial data
- 2 **SOMbrero**





- **SOMbrero** is an R package implementing stochastic variants of SOM for non vectorial data (see **yasomi** for batch versions)
- first release: March 2013; latest release: November 2013 (version 0.4-1)
- depends on R (version ≥ 3.0) <http://www.r-project.org>



and on several packages available on CRAN:

**wordcloud, igraph, RColorBrewer, scatterplot3d, knitr,
shiny**

- available at <http://sombbrero.r-forge.r-project.org> (licence GPL) and can be installed from inside R using

```
install.packages("SOMbrero",  
  repos="http://R-Forge.R-project.org")
```



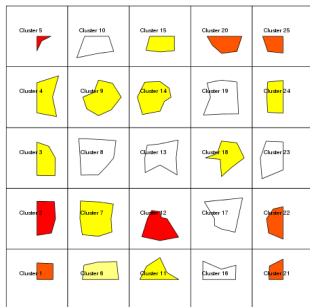
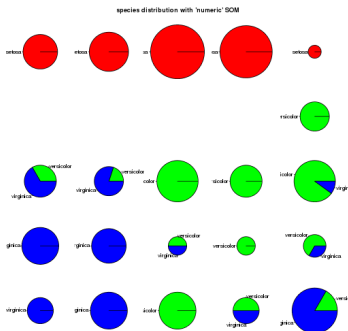


- 1 **3 algorithms** available through one function `trainSOM`
 - numeric SOM (input: $(n \times p)$ -matrix with n observations of p variables)
 - KORRESP (input: $(p \times q)$ -contingency table)
 - relational SOM (input: $(n \times n)$ -dissimilarity matrix for n individuals)



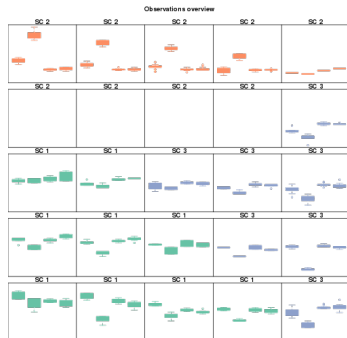
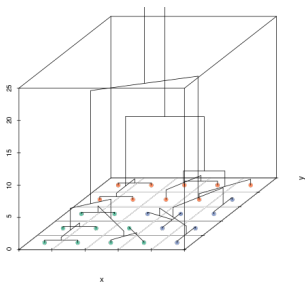


- 1 **3 algorithms** available through one function **trainSOM**
- 2 **many graphics** available through one function **plot** with two main arguments **what** (prototypes, observations, additional variable) and **type** (color, 3d, barplot, poly.dist, words, pie...)





- 1 3 algorithms available through one function `trainSOM`
- 2 many graphics
- 3 super-clustering (HC on prototypes) with associated graphics through functions `superClass` and `plot`





- 1 3 algorithms available through one function `trainSOM`
- 2 many graphics
- 3 super-clustering (HC on prototypes) with associated graphics
- 4 quality measures (quantization error, topographic error) with quality



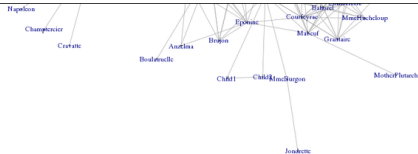


- **3 datasets** corresponding to the three algorithms (**iris**, **presidentielles2002** and **lesmis**, a graph from « Les Misérables »)





- **3 datasets** corresponding to the three algorithms (iris, presidentielles2002 and lesmis, a graph from « Les Misérables »)
- comprehensive (HTML) **vignettes** included in the package and available on the website



The `dissim.lesmis` object is a matrix with entries equal to the length of the shortest path between two characters (obtained with the function `shortest.paths` of `pac` characters' names to ease the use of the graphical functions of `SOMbrero`).

Training the SOM

```
set.seed(4031719)
mis.som = trainSOM(x.data = dissim.lesmis, type = "relational", nb.save = 10,
  init.proto = "random")
plot(mis.som, what = "energy")
```

Energy evolution





- **3 datasets** corresponding to the three algorithms (**iris**, **presidentielles2002** and **lesmis**, a graph from « Les Misérables »)
- comprehensive (HTML) **vignettes** included in the package and available on the website
- **Web User Interface** (made with **shiny**) for using the package even if you do not know R programming language (included in the package or available online at <http://shiny.nathalievilla.org/sombrero> but can be very slow)

SOMbrero Web User Interface (v0.1)

Select the data type:
Numeric

Import Data Self-Organize Plot Map Superclasses Combine with external information Help

Third step: plot the self-organizing map

In this panel and the next ones you can visualize the computed self-organizing map. This panel contains the standard plots used to analyze the map.

Options

Plot what?
Prototypes

Type of plot:
polygon distances

Show cluster names

Welcome to SOMbrero, the open-source on-line interface for self-organizing maps (SOM).
This interface trains SOM for numerical data, contingency



Let's go into **SOMbrero**...



disclaimer: as is standard during demos, something nasty might happen and nothing would work due to some weird technical issues... and the speaker would look like an idiot!





SOMbrero

- is easy to use (with a simple graphical interface)
- can be used with various data
- contains many tools for interpreting the results





SOMbrero

- is easy to use (with a simple graphical interface)
- can be used with various data
- contains many tools for interpreting the results
- ... has unfortunately been implemented by girls so default colors may not be suited for men (but they can easily change them)

Perspectives

- speed up the code
- more quality criteria
- more options (i.e., Gaussian neighborhood, weighted observations...)





Thank you for your attention...



... questions?





Cottrell, M., Letrémy, P., and Roy, E. (1993).

Analyzing a contingency table with Kohonen maps: a factorial correspondence analysis.

In Cabestany, J., Mary, J., and Prieto, A. E., editors, *Proceedings of International Workshop on Artificial Neural Networks (IWANN 93)*, Lecture Notes in Computer Science, pages 305–311. Springer Verlag.



Goldfarb, L. (1984).

A unified approach to pattern recognition.

Pattern Recognition, 17(5):575–582.



Hammer, B. and Hasenfuss, A. (2010).

Topographic mapping of large dissimilarity data sets.

Neural Computation, 22(9):2229–2284.



Kohonen, T. (2001).

Self-Organizing Maps, 3rd Edition, volume 30.

Springer, Berlin, Heidelberg, New York.



Olteanu, M. and Villa-Vialaneix, N. (2014).

On-line relational and multiple relational SOM.

Neurocomputing.

Forthcoming.

