

Three-way clustering around latent variables approach with constraints to improve configurations' interpretability

Véronique Cariou^a, Tom F. Wilderjans^{b,c}

^a StatSC, ONIRIS, INRA, 44322 Nantes, France

^b Methodology and Statistics Research Unit, Leiden University, The Netherlands

^c Research Group of Quantitative Psychology and Individual Differences, Leuven, Belgium

XXVIe Rencontres de la Société Francophone de Classification
3 – 5 sept 2019, Nancy



Introduction

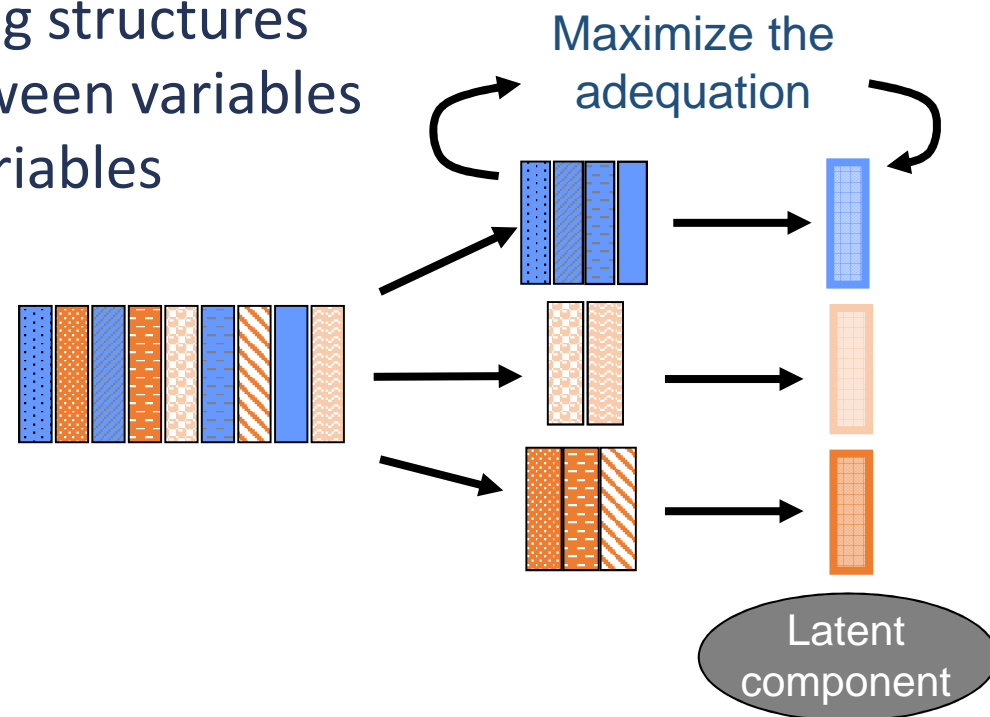
Clustering of variables around latent components with CLV:

- ✓ Understand the underlying structures
- ✓ Detect redundancies between variables
- ✓ Reduce the number of variables

Maximize

$$T^{(Q)} = \sum_{q=1}^Q \sum_{j=1}^J \delta_{qj} \text{cov}^2(\mathbf{x}_j, \mathbf{c}_q)$$

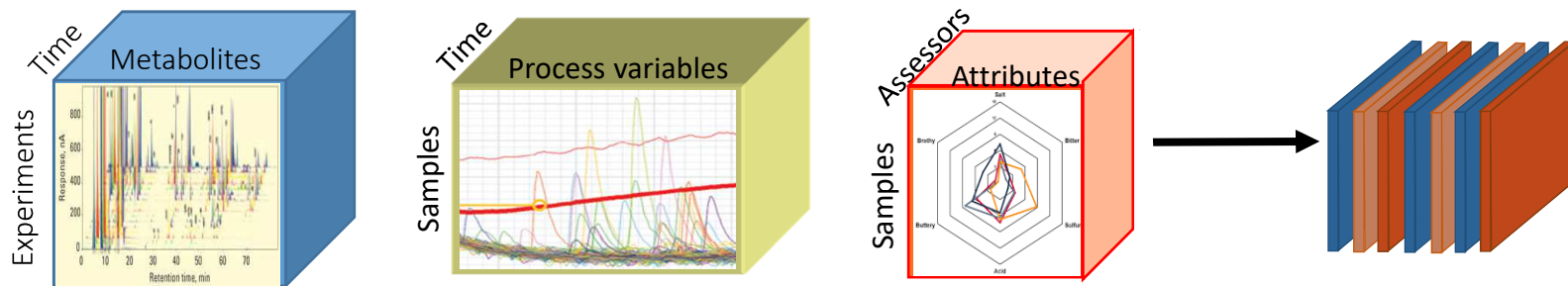
with $\text{var}(\mathbf{c}_q) = 1$



Vigneau, E., & Qannari, E. M. (2003). Clustering of variables around latent component. *Communications in Statistics – Simulation and Computation*, **32**,1131–1150

Introduction

Extension of CLV to three-way arrays with CLV3W:



Connexion with clusterwise models and other clustering techniques

Basford, K. E., et McLachlan, G. J. 1985. The mixture method of clustering applied to three-way data. *Journal of Classification*, **2**(1), 109–125

Vichi, M. 1999. One-mode classification of a three-way data matrix. *Journal of Classification*, **16**(1), 27–44

Llobell, F., Cariou, V., Vigneau, E., Labenne, A., et Qannari, E. M. 2018. Analysis and clustering of multiblock datasets by means of the stasis and clustasis methods. application to sensometrics. *Food Quality and Preference*.

Overview

Clustering around latent components for three-way data

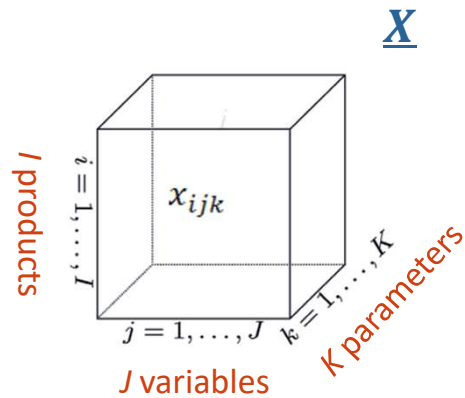
Adding a non negativity constraint on loadings

Application on sensory data for consumers' segmentation

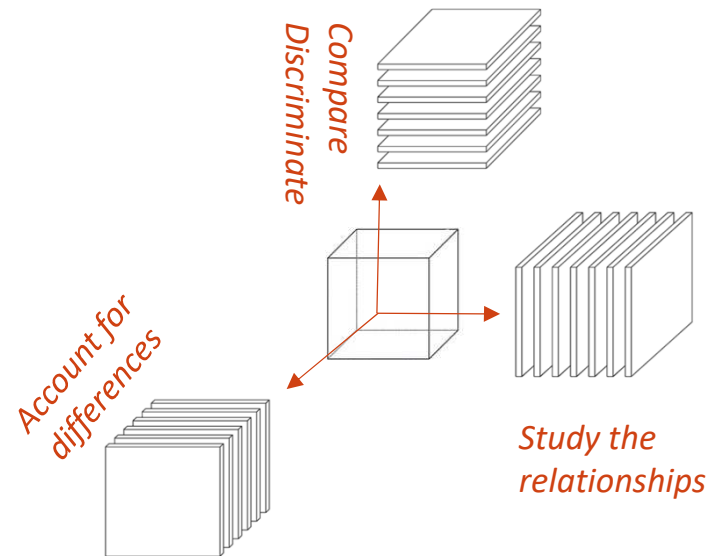
Imposing a global weighting scheme

Conclusion and perspectives

Clustering around latent components for Three-way data



x_{ijk} : score associated with the i^{th} product rated by the j^{th} consumer along with the k^{th} parameters



\underline{X}_j – j^{th} lateral slice of \underline{X} – contains the scores of I samples on the j^{th} variable according to K parameters.

Clustering around latent components for Three-way data

Three-way partitioning of the 2nd mode:

- ✓ Maximise the adequation between each lateral slice and the estimated one associated with its group

Minimize the Loss function:

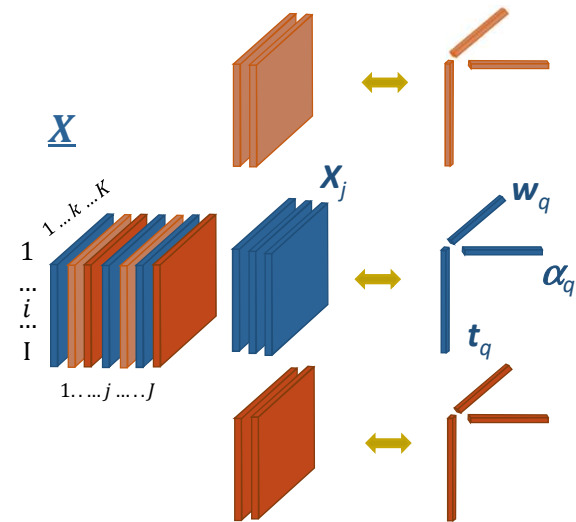
$$f = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \|\mathbf{X}_j - \alpha_{jq} (\mathbf{t}_q \mathbf{w}_q^T)\|_F^2$$

1st mode latent component: \mathbf{t}_q

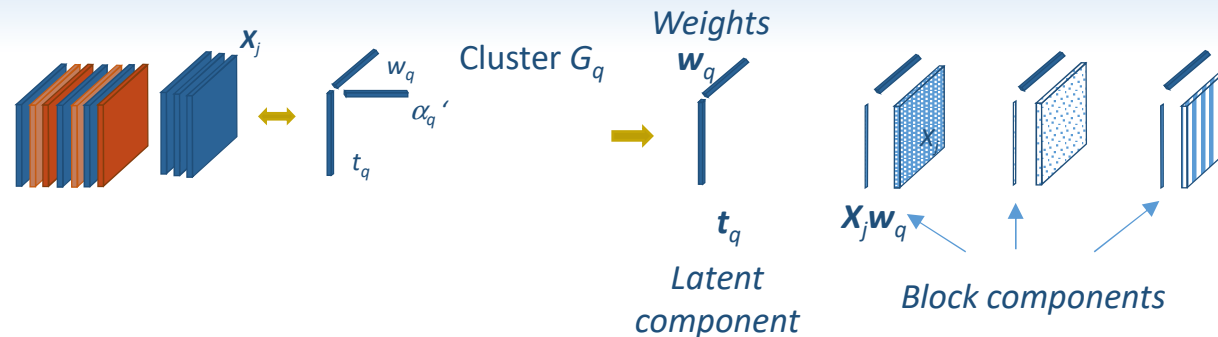
2nd mode loadings of variables: α_q

3rd mode attributes' weights: \mathbf{w}_q

cluster dummy: δ_{jq}



Clustering around latent components for Three-way data



Minimize the Loss function:

$$f = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \|\mathbf{X}_j - \alpha_{jq} (\mathbf{t}_q \mathbf{w}_q^T)\|_F^2$$

↔ Maximise the squared covariance function:

$$g = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \text{cov}^2(\mathbf{X}_j \mathbf{w}_q, \mathbf{t}_q)$$

with $\|\mathbf{w}_q\|=1$ and $\|\mathbf{t}_q\|=1$

Krijnen, W. P. 1993. The analysis of three-way arrays by constrained PARAFAC methods. DSWO Press
 Wilderjans T. F., et Ceulemans, E. 2013. Clusterwise Parafac to identify heterogeneity in three-way data.
Chemometrics and Intelligent Laboratory Systems, **129**, 87–97

Initialisation

- Using a multi-start procedure obtained by randomly assigning the variables to Q clusters
- Using a rational initial partitioning obtained by applying a Hierarchical Ascendant Algorithm based on criterion f and Ward's aggregation criterion
 - The variation of the criterion from step l to step $l+1$ corresponds to the aggregation of two clusters, say G_A and G_B .
 - It can be written as:

$$\Delta f = - \sum_{X_j \in G_A} \|X_j - \alpha_{jG_A}(\mathbf{t}_{G_A} \mathbf{w}_{G_A}^\top)\|_F^2 - \sum_{X_j \in G_B} \|X_j - \alpha_{jG_B}(\mathbf{t}_{G_B} \mathbf{w}_{G_B}^\top)\|_F^2 \\ + \sum_{X_j \in (G_A \cup G_B)} \|X_j - \alpha_{j(G_A \cup G_B)}[\mathbf{t}_{(G_A \cup G_B)} \mathbf{w}_{(G_A \cup G_B)}^\top]\|_F^2$$

Algorithm

Starting from an initial partitioning into Q clusters

Iterate the following two steps until convergence:

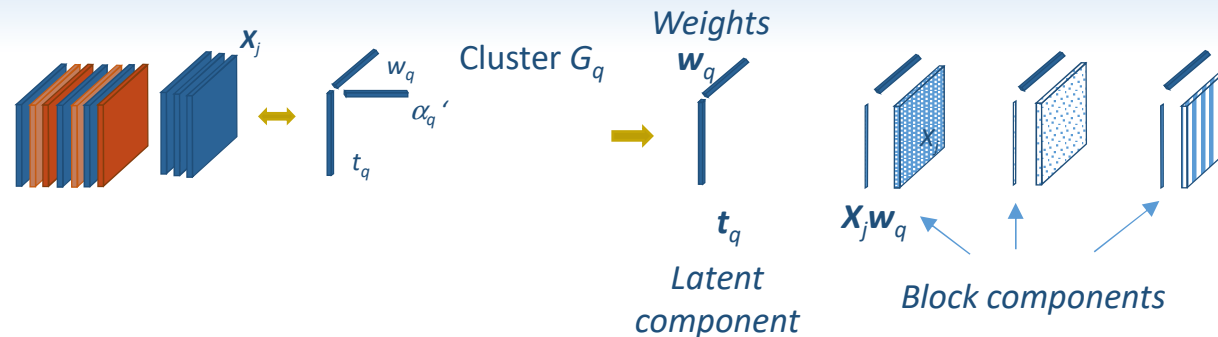
- (1) computing the cluster-specific parameters conditional on the cluster memberships
- (2) updating the cluster membership δ_{jq} of each variable conditional on the cluster-specific parameters (i.e., \mathbf{t}_q and \mathbf{w}_q)

For the assignment, variable j is assigned to the cluster G_q for which

f_{jq} is minimal: $f_{jq} = \|\mathbf{X}_j - \alpha_{jq}(\mathbf{t}_q \mathbf{w}_q^\top)\|_F^2$ with $\|\mathbf{t}_q\|=1, \|\mathbf{w}_q\|=1$

Determination of α_{jq} by means of a linear regression

Adding a non negativity constraint



Minimize the Loss function:

$$f = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \|X_j - \alpha_{jq} (t_q w_q^T)\|_F^2$$

with $\|w_q\|=1$, $\alpha_{jq} \geq 0$ and $\|t_q\|=1$

↔ Maximise the covariance function:

$$g = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \text{cov}^2(X_j w_q, t_q)$$

Determination of α_{jq} by means of a non negative linear regression

Consumers' segmentation

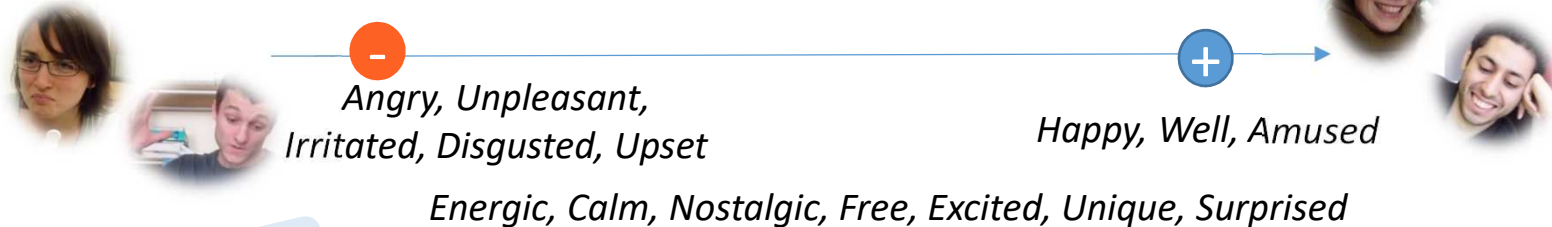
Emotion ratings associated to coffee aromas

- ✓ 84 undergraduate students from Oniris
- ✓ 12 coffee aromas

Earth, Hay, Cedar; Vanilla, Coriander seeds, Flower Coffee, Apricot, Lemon, Honey, Basmati Rice, Hazelnut, Medicine



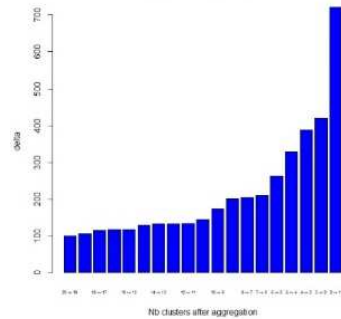
- ✓ 15 emotions rated with a 5-point Likert scale



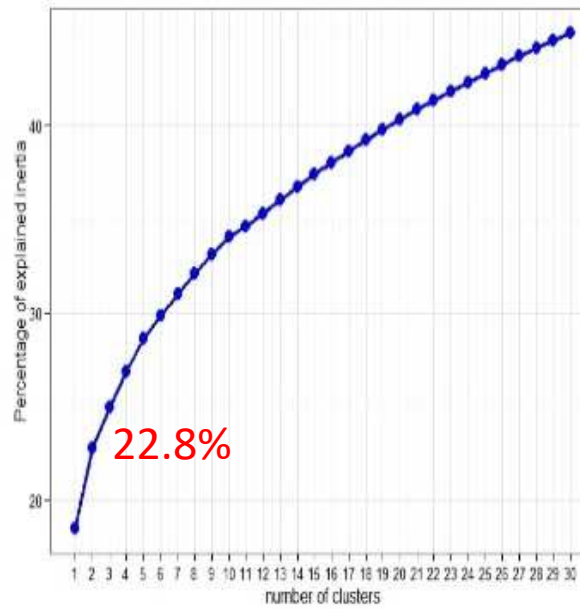
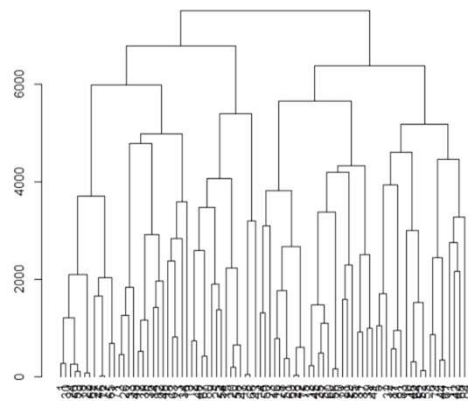
within a cluster

consumer agreement on the product ratings

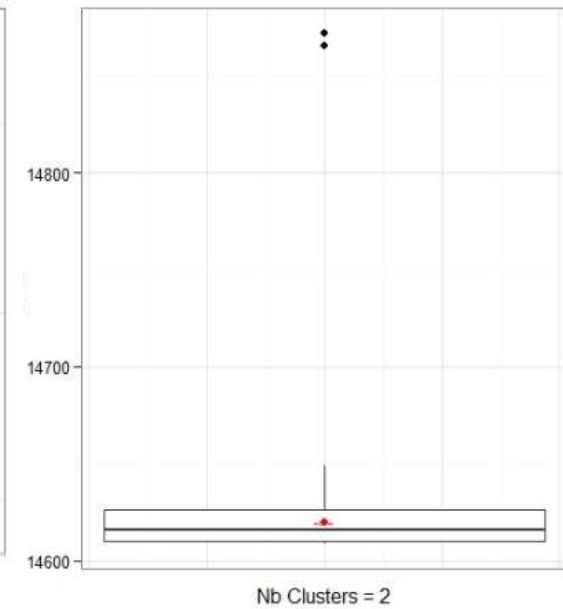
Consumers' segmentation



Evolution of the aggregation criterion



% of explained inertia / number of clusters



Sensitivity to initialization

Consumers' segmentation

Functions in Package ClustVarLV

```
res <- CLV3W_kmeans(coffee,K=2,mode.scale=2, NN=TRUE,init=50,cp.rand=5)
summary(res)
```

```
$size
clusters
 1  2
42 42

$prop_explained_per_cluster
[1] 21.308 24.316

$prop_explained_total
[1] 22.812

$comp
      Comp1  Comp2
vanilla -0.306  0.061
B.Rice  -0.058 -0.218
Lemon   0.524  0.435
Coffee.Flower 0.422  0.278

$weight
      Comp1  Comp2
calm      0.215  0.179
nostalgic 0.188  0.227
angry     -0.224 -0.222
disgusted -0.377 -0.371
unique     0.094  0.127
excited   0.138  0.207
unpleasant -0.337 -0.327
disappointed -0.232 -0.281
free      0.229  0.188
surprised -0.008 -0.009

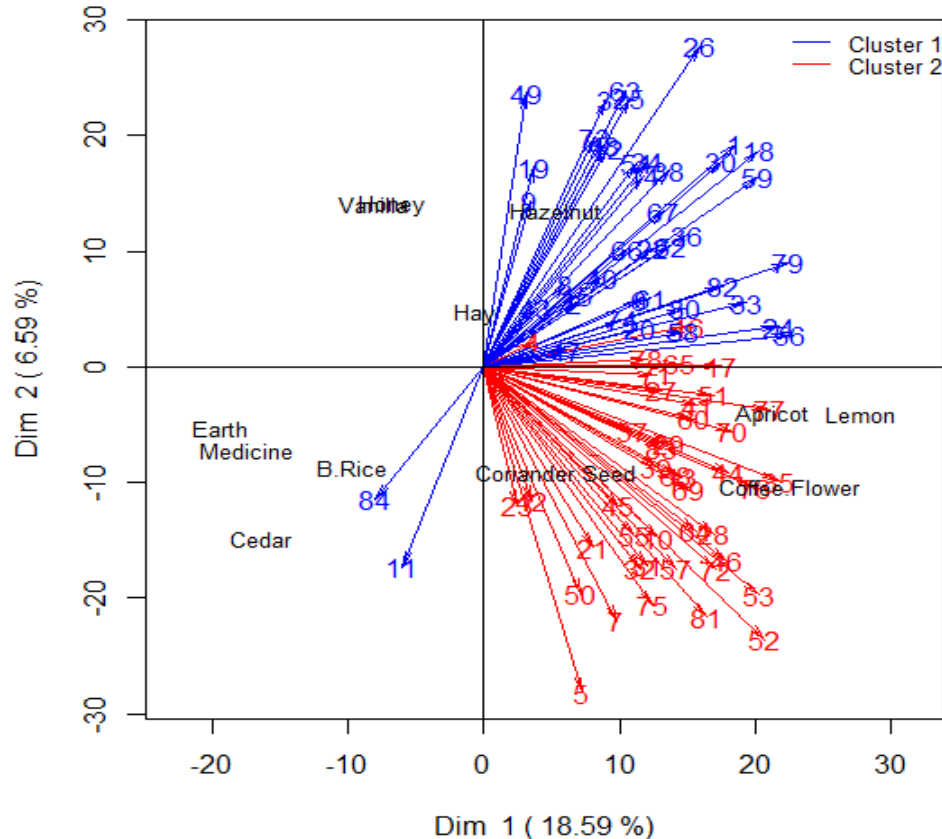
$scormatrix
      Comp.1  Comp.2
Comp.1 1.000  0.586
Comp.2 0.586  1.000

$groups
$groups$`cluster 1`
  loading cor in.group cor next.group
1  0.218   0.824   0.504
2  0.078   0.605   0.343
3  0.152   0.600   0.135
6  0.134   0.623   0.408
8  0.086   0.384   0.157
9  0.072   0.528  -0.055
11 0.000  -0.434  -0.115
12 0.128   0.654   0.226
14 0.154   0.692   0.345
15 0.090   0.366   0.160
18 0.236   0.808   0.505
19 0.091   0.376  -0.043
20 0.129   0.536   0.402
22 0.144   0.639   0.382
```

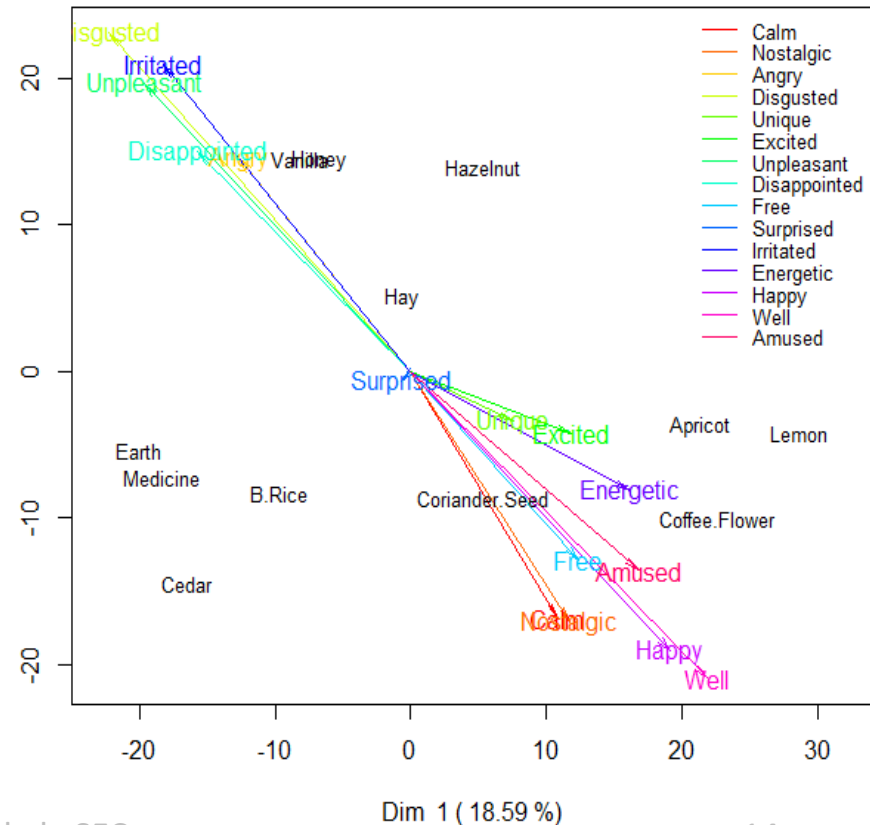
Consumers' segmentation

```
plot_var.clv3w(res.clv3w,K=2,labels=TRUE,cex.lab=0.8,beside=FALSE,mode3=TRUE)
```

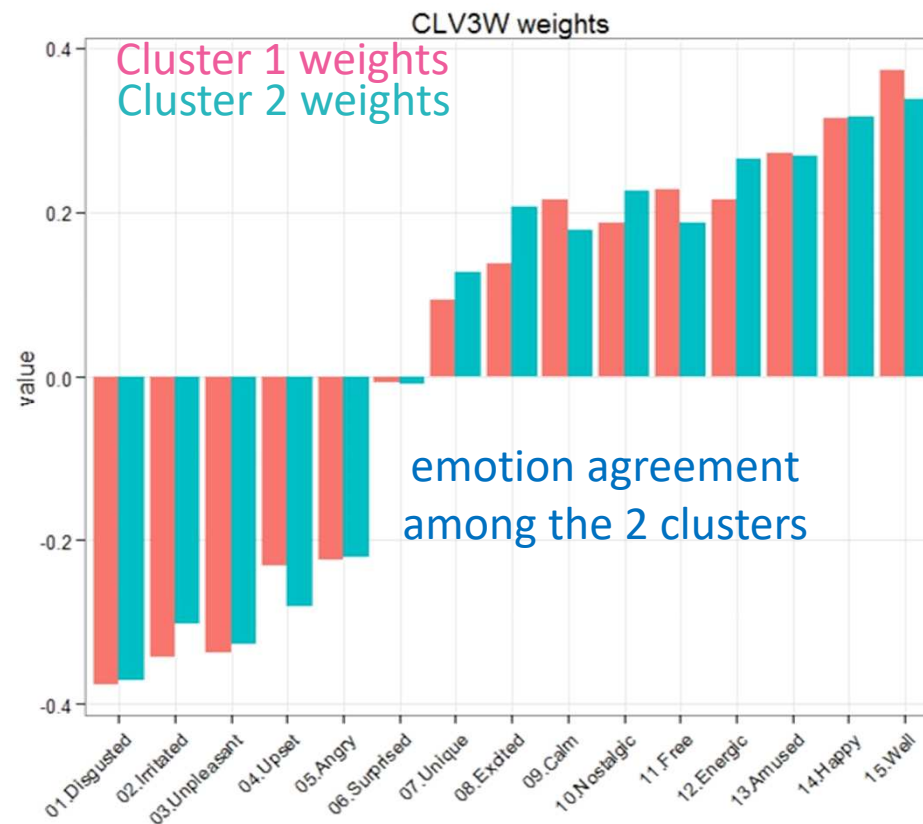
Candecomp Parafac Scores plot



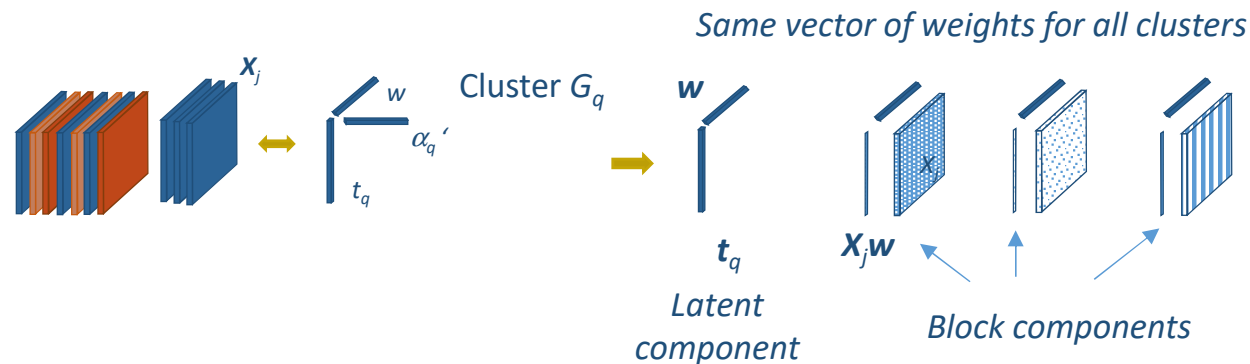
Candecomp Parafac Scores plot



Consumers' segmentation



Imposing a global weighting scheme



Minimize the Loss function:

$$f = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \|\mathbf{X}_j - \alpha_{jq} (\mathbf{t}_q \mathbf{w}^T)\|_F^2$$

↔ Maximise the squared covariance function:

$$g = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \text{cov}^2(\mathbf{X}_j \mathbf{w}, \mathbf{t}_q)$$

with $\|\mathbf{w}\|=1$ and $\|\mathbf{t}_q\|=1$

Algorithm

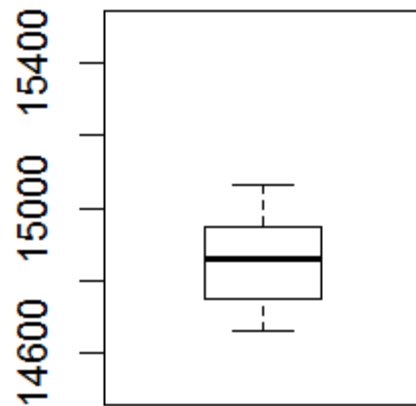
Iterate the following two steps until convergence:

(1) computing parameters conditional on the cluster memberships:

- initializing \mathbf{w} , the two steps are alternated until convergence:
- Each \mathbf{t}_q is updated as the left singular vector corresponding to the largest singular value of the $(I \times \#G_q)$ matrix of the block components associated with the variables of G_q : $[\cdots |\mathbf{X}_j \mathbf{w}| \cdots]_{\mathbf{X}_j \in G_q}$
- Global \mathbf{w} is updated as the left singular vector associated with the largest singular value from of the $(K \times J)$ matrix $[\cdots |\mathbf{X}_j^\top \mathbf{t}_q| \cdots]_{j=1, \dots, J}$

(2) updating the cluster membership δ_{jq} of each variable conditional on the cluster-specific \mathbf{t}_q and the global weighting scheme \mathbf{w}

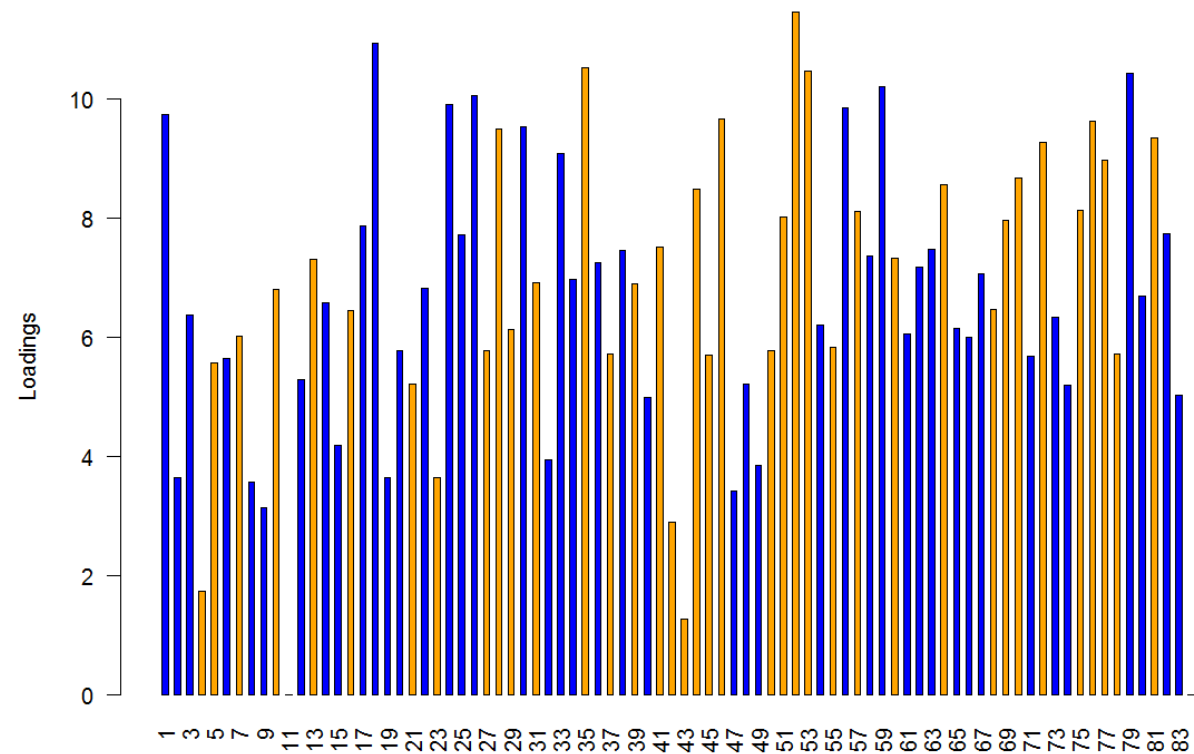
Consumers' segmentation



Sensitivity to initialization

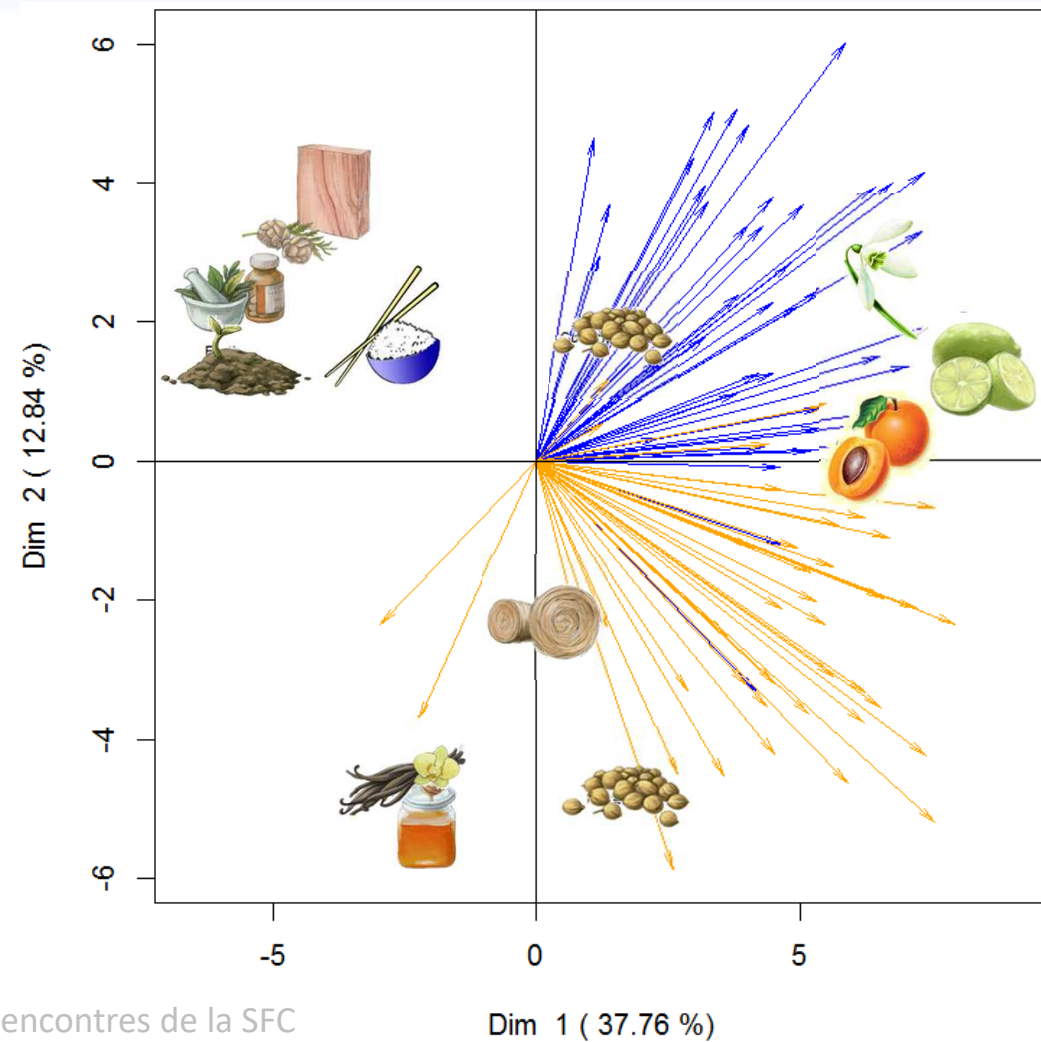
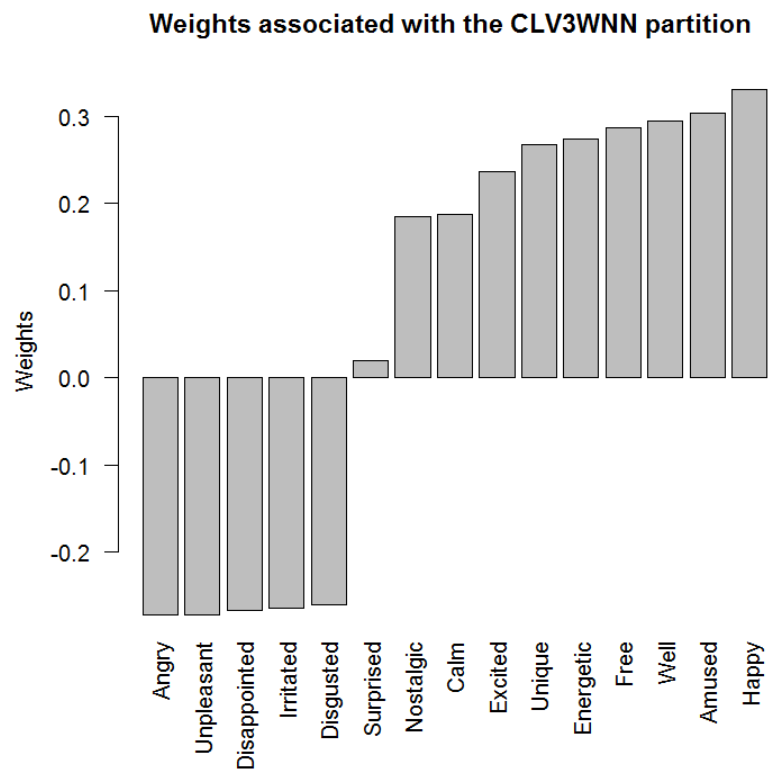
Inertia explained with NN
constraint: 22.8%

Inertia explained with NN
and Global Weights: 22.5%



Size of cluster 1: 44 ; Size of cluster 2: 38
2 consumers set aside

Consumers' segmentation



Conclusion

Cluster analysis of three-way data with CLV3W providing:

- ✓ Latent component associated with the 1st mode samples
- ✓ Loadings reflecting the adequation of each slice in its cluster
- ✓ Weights deriving a block component for each slice

Available in Cran R

Package 'ClustVarLV'

June 28, 2019

Title Clustering of Variables Around Latent Variables

Version 2.0.0

Author Evelyne Vigneau [aut, cre], Mingkun Chen [ctb], Veronique Cariou [aut]



Easier interpretation with constraints on loadings and weights

Perspectives

Choice of the number of clusters:

✓ Adaptation of Hartigan criterion $H_Q = (J - Q - 1) \left(\frac{f_Q}{f_{Q+1}} - 1 \right)$

Index corresponding to the largest difference between two consecutive ones

Choice between full model or common weights constraint

Further comparisons needed with other strategies such as simultaneous decomposition and clustering models

Minimize the Loss function: $f = \sum_{j=1}^J \sum_{q=1}^Q \delta_{jq} \| \mathbf{X}_j - \mathbf{t}_q \mathbf{w}_q^\top \|^2_F$ with $\| \mathbf{w}_q \| = 1$

Vichi, M., Rocci, R., et Kiers, H. A. 2007. Simultaneous component and clustering models for three-way data : within and between approaches. *Journal of Classification*, **24**(1), 71–98

Thank you for your attention

