

| | | | | | | | |
|---------|-------------------|----------------------------|---|------------|---------------|------------|-----|
| Densité | Clustering ○○○ | Clustering. Images ○○○○ | ACP et la famille ○○○○○○ ○○○○○○○○○○ ○○○○○○ ○○○○○○ | ICA ○○○ | MDS ○○○○○○ | Ressembler | TP3 |
|---------|-------------------|----------------------------|---|------------|---------------|------------|-----|

Apprentissage Non Supervise

B. Ghattas & G. Oppenheim

Université d'Aix-Marseille - Université Paris-Sud

December 8, 2016

| | | | | | | | |
|---------|-------------------|----------------------------|--|------------|---------------|------------|-----|
| Densité | Clustering ○○○ | Clustering. Images ○○○○ | ACP et la famille ○○○○○○ ○○○○○○○○ ○○○○○○ ○○○○○ | ICA ○○○ | MDS ○○○○○○ | Ressembler | TP3 |
|---------|-------------------|----------------------------|--|------------|---------------|------------|-----|

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

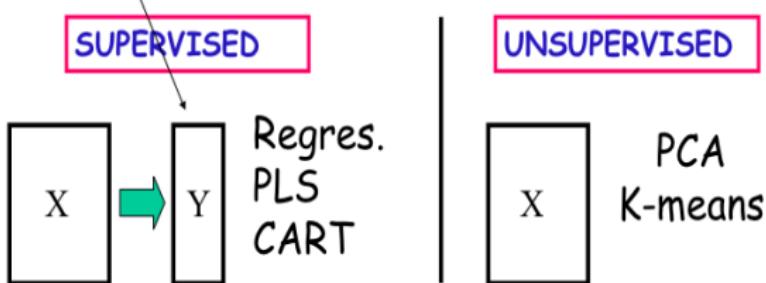
7 Ressembler

8 TP3

To be or not to be Supervised?

Psychological studies: when there is an information coming from a SUPERVISOR who says: It's good or It isn't.

Here the context is a little different:
there is a target



The last application I studied is a Supervised problem. Quality control of a Rubber/Metal suspension bushes for an engine mount
(articulation élastique pour support moteur)

7

2/101



Les méthodes de cette séance

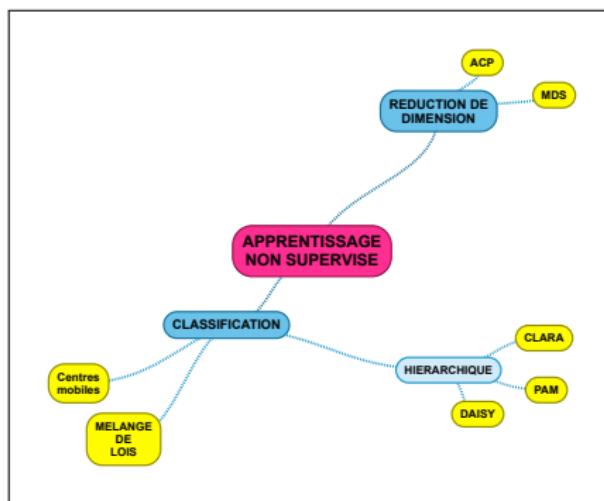


Figure: Les méthodes Non-Supervisées de cette séance



Exemple détaillé

Le cas étudié est un problème de cytométrie consistant à analyser des particules contenues dans un échantillon prélevés dans un environ marin.



An application: Flow cytometer dataset

- A collection of Phytoplankton cells was performed in a sailing yacht harbour in Marseille.
- Samples were collected every 30 min during one day, from 26-09-2011 at 17h00 to 27-09-2011 at 17h00.
- Five signals were recorded for each cell: FWS (ForWard Scatter), SWS (Side-Ward Scatter), FLY (Fluorescence Yellow), FLO (Fluorescence Orange) and FLR (Fluorescence Red).
- From each signal the following features were computed: average, maximum, length, inertia, fill factor, center of gravity and asymmetry.

External factors

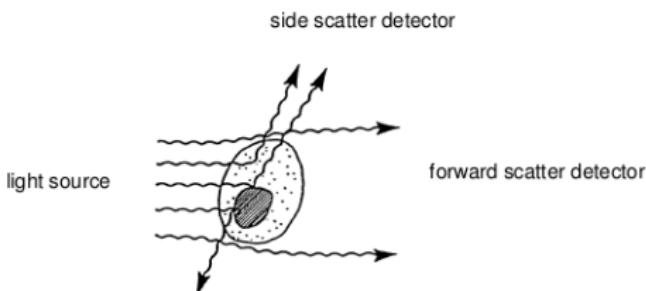
Some meteorological features were collected during those 2 days (26-27 September 2011) in order to characterize the evolution of Phytoplankton particles:

- - Temperature of the seawater (in Celsius °C).
- - Acidity/Basicity of the seawater (in pH units).
- - Oxidation-Reduction Potential (ORP): a measure of the cleanliness of the seawater and its ability to break down contaminants. (in mV).
- - Salinity of the water (in ppt \$parts per thousand).
- - Turbidity of the seawater: measures the cloudiness or haziness of the water. (in Nephelometric Turbidity Units (NTU)).
- - Chlorophyll (in ?g per liter ?g/l): a molecule that absorbs sunlight and uses its energy to synthesize carbohydrates from CO₂ and water, its concentration in the seawater gives an insight of the distribution of Phytoplankton.



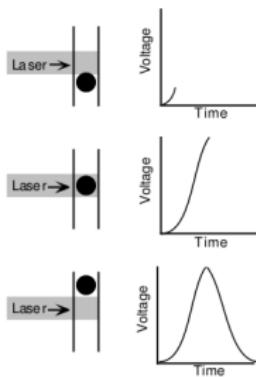
Data acquisition

Detection of the signal in the Flow Cytometer: Light scattering occurs when a particle deflects incident laser light. A detector in front of the laser beam measures forward scatter FS and several detectors to the side measure side scatter SS. Fluorescence detectors measure the fluorescence emitted from stained cells.



Signal creation

The highest point of the signal occurs when the particle is in the center of the beam and the maximum amount of scatter or fluorescence is achieved.



The Forward Scatter light FWS is proportional to the size of the cell and the Side-ward Scatter light is proportional to the internal complexity of the cell. The three Fluorescences (Yellow, Orange and Red) are used for the separation of the particles by a particular protein.



More about the dataset..

We have 39 samples, each consisting on almost 100.000 observations, with 50 variables corresponding to the features extracted from the signals.

Objective

- Cluster each sample. Choose an appropriate method, and an appropriate number of clusters.
- How can we model the differences between successive partitions ?
- How can we relate partitions' evolution to external factors like Temperature, turbidity or salinity ? Which ones are relevant to explain this evolution ?

Densité

Clustering
○○○

Clustering. Images
○○○○

ACP et la famille



ICA



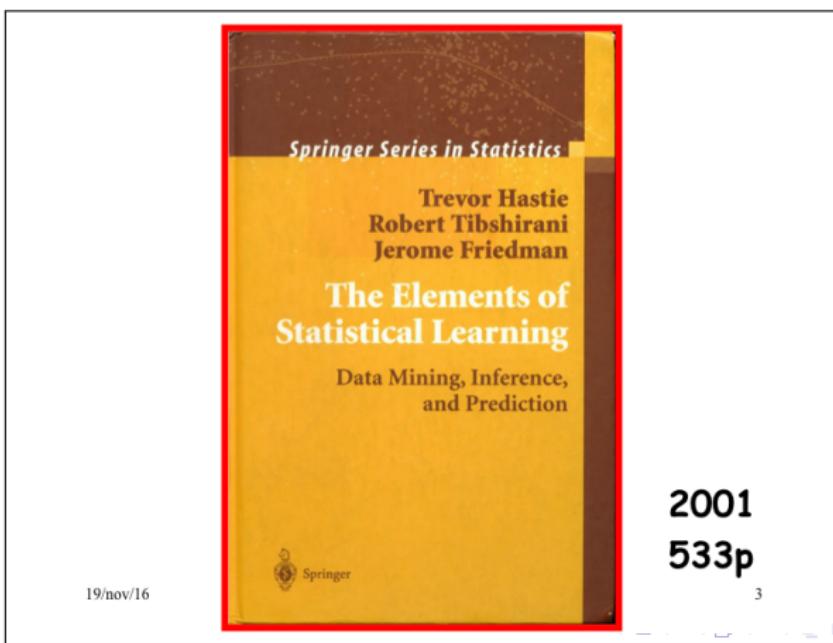
MDS



Ressembler

TP3

Référence





le cours est appuyé sur

- Le MOOC de Hastie sur le non-supervisé
- Le livre de Hastie-Tibshirani-Friedmann

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

Densité

Déterminer la densité.

- *Histogramme*
- *Densité: noyaux*

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

Clustering methods

Clustering aims to partition the data. Partitioning may be:

- *Hierarchical*: find successive groups splitting or joining previously established groups ("bottom-up", "top-down")
- *Non Hierarchical*: k-means. Density based methods: Mixture models, DBSCAN.

Some **supervised** methods are based also on partitioning the space of explanatory variables (CART, SVM, LDA).

Clustering methods

- Works often with **dissimilarity** matrix.
- Often needs the number k of clusters to be fixed in advance.
- There exist some heuristics to estimate k (pseudo T, pseudo F).
- We need also to estimate the quality of a partition and to compare partitions.

Often we need to deal with ...

- qualitative datasets, numeric, or mixed
- functional data, image data

The most widely used method: The k -means algorithm.

k-means

- 1** Choose arbitrary centers,
- 2** Repeat the two following steps until the centers do not change.
 - 1** assign each observation to the closest center,
 - 2** compute the new center of each class using the observations assigned to that class,
- 3** Repeat the previous steps 10 times and keep the partition $(G_{jk}, j = 1\dots)$ with the minimum within-cluster sum of squares given by,

$$\sum_{i=1}^n \sum_{j=1}^k \| \mathbf{x}_i - c_j \|^2 \mathbf{1}_{\{i \in G_j\}}$$

where G_j is the j -th group and c_j is the corresponding center.

- 4** Create SRONG FORMS.



Kmeans

Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

```

○○○○○○○
○○○○○○○○
○○○○○○○
○○○○
○○○○○

```

ICA

○○○

MDS

○○○○○○○

Ressembler

TP3

HastieKmeans1

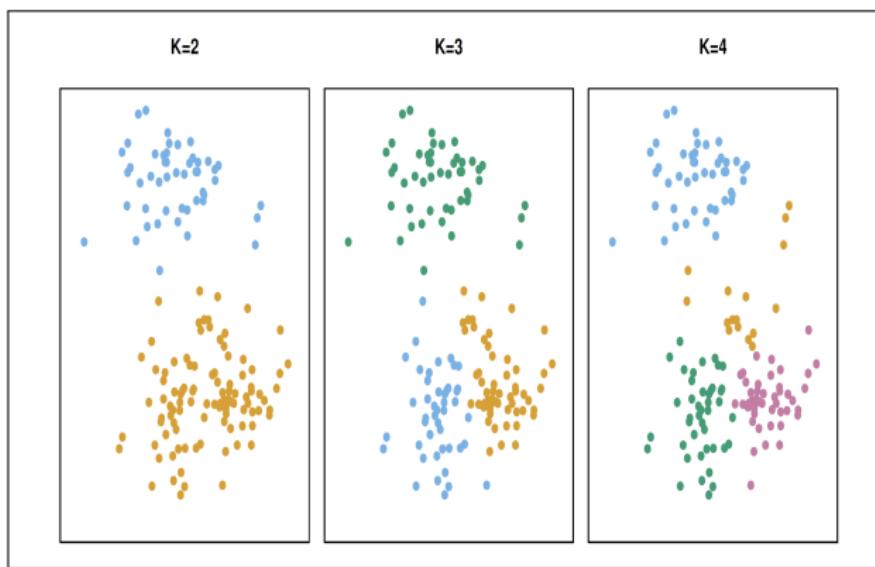


Figure: HastieKmeans1

Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

```

○○○○○○○
○○○○○○○○
○○○○○○○
○○○○
○○○○○

```

ICA

○○○

MDS

○○○○○○○

Ressembler

TP3

HastieKmeans2

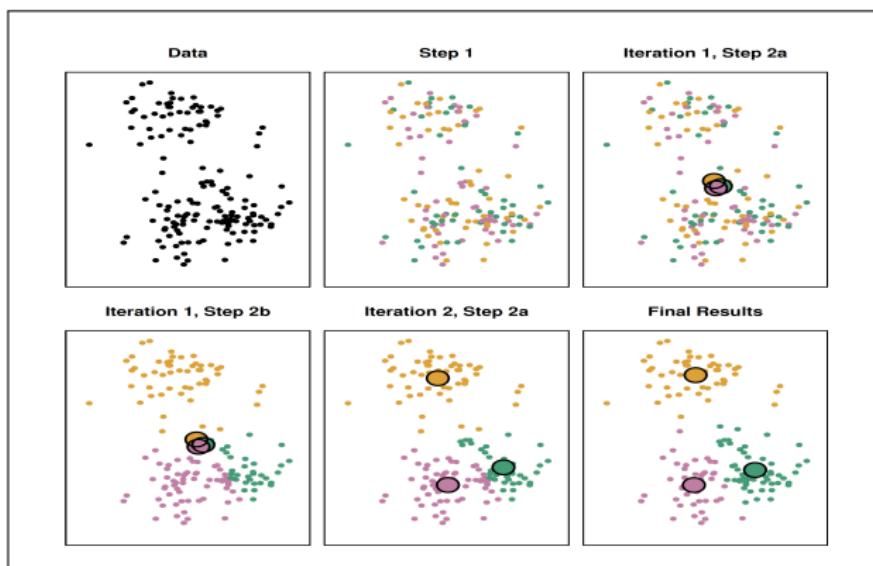


Figure: HastieKmeans2

Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

```

○○○○○○○
○○○○○○○○
○○○○○○○○
○○○○○
○○○○○○

```

ICA

○○○

MDS

○○○○○○○

Ressembler

TP3

HastieKmeans3

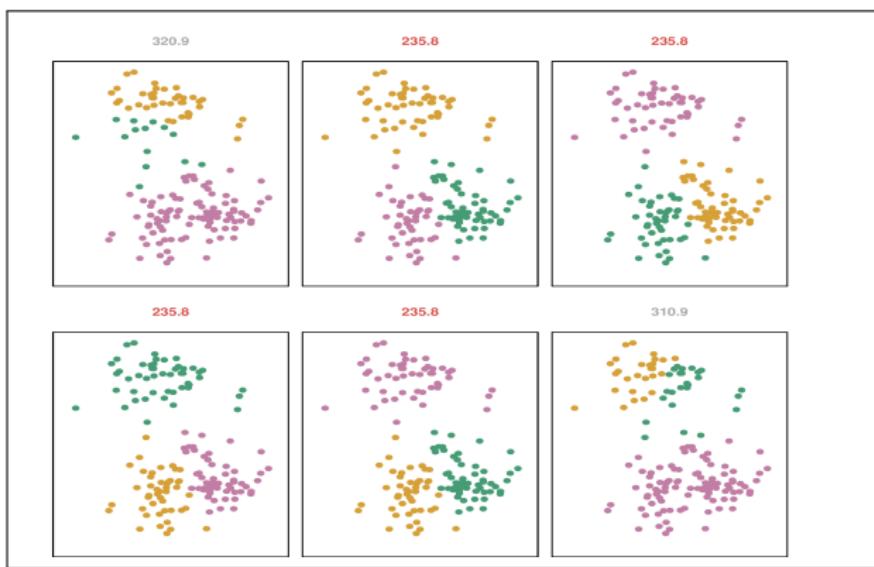


Figure: HastieKmeans3

k-modes/*k*-medians

Extensions of the well-known *k*-means algorithm [8].

- The *k*-modes algorithm [3] for categorical data that extends the *k*-means algorithm. It seeks to partition the observations into *k* groups such that the distance from the observations to the cluster modes is minimized. Simple-matching distance to assess the dissimilarity between pairs of observations is often used (Number of disagreements between pairs/Total number of pairs).

See the *klaR* package [11] from R.

- The *k*-medians approach is recommended in [6] for addressing ordinal data. It is similar to the *k*-means algorithm except that it uses medians rather than means as centers for the clusters. Manhattan distance ($\sum_i |x_i - y_i|$) may be used to assess the dissimilarity between the observations and the cluster medians.

See the *flexclust* package [6] from R.



Hierarchical Clustering

Hierarchical methods [9] aim to construct a dendrogram, given a distance d , and they can be

- 1** ascendant ("bottom-top") or
- 2** descendant ("top-down").

In ascendant hierarchical clustering, the hierarchy is constructed by iterative aggregation of pairs of nearest clusters. The distance between two clusters may be defined in several ways, the most popular being

- simple
- average
- complete linkage.

For ordinal data one may use the Manhattan distance as the dissimilarity measure between the observations, or the mutual information for nominal data.

The algorithm is implemented in the *hclust* function from R.

Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

○○○○○○
○○○○○○○○
○○○○○○○○○
○○○○○○
○○○○○○○○

ICA

○○○

MDS

○○○○○○○

Ressembler

TP3

PAM, CLARA , DAISY

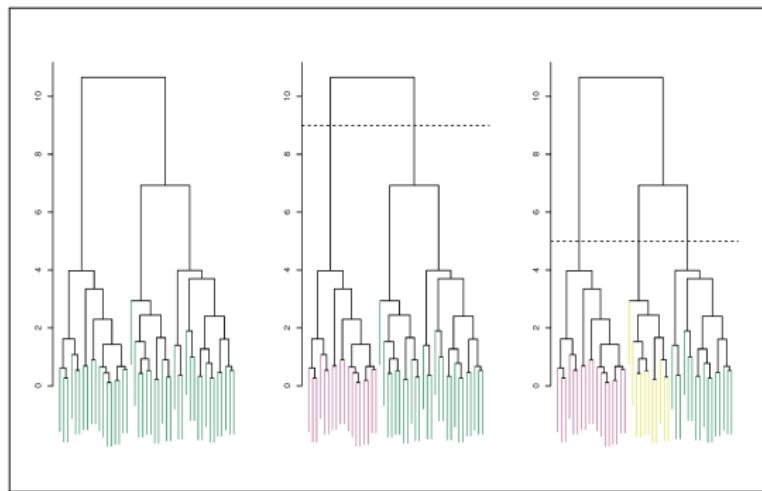


Figure: HastieCluster1

HastieCluster2

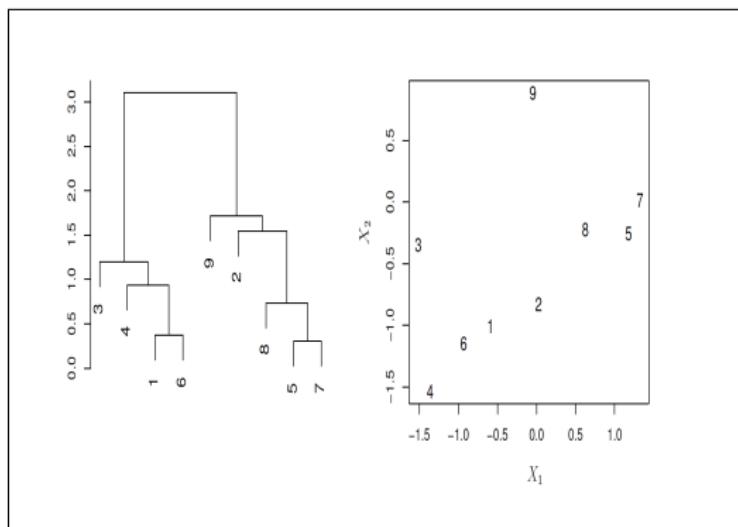


Figure: HastieCluster2

HastieCluster3

Types of Linkage

| <i>Linkage</i> | <i>Description</i> |
|----------------|---|
| Complete | Maximal inter-cluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>largest</i> of these dissimilarities. |
| Single | Minimal inter-cluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>smallest</i> of these dissimilarities. |
| Average | Mean inter-cluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>average</i> of these dissimilarities. |
| Centroid | Dissimilarity between the centroid for cluster A (a mean vector of length p) and the centroid for cluster B. Centroid linkage can result in undesirable <i>inversions</i> . |

Figure: HastieCluster3



Category Utility ([3])

Consider a partition, $C = C_k$, $k = 1, \dots, K$, found by a clustering algorithm based on given Variables X_j , $j = 1, \dots, p$. The attributes are assumed nominal so that each X_j has a set of values or categories, x_{jl} . The category utility function scores partition C against the set of variables according to formula:

$$CU(C) = \frac{1}{K} \sum_{k=1}^K P(C_k) \left[\sum_j \sum_l P(X_j = x_{jl} | C_k)^2 - \sum_j \sum_l P(X_j = x_{jl})^2 \right]$$

Probability matching is random decision strategy using the empirical distribution, which is different from the bayesian majority vote prediction. The probability of winning when distribution is given by p_j , is $\sum_j p_j^2$



The matching error

We assess the performance using the matching error.

Let y_1, \dots, y_n be the class labels of each observation, and let $\hat{y}_1, \dots, \hat{y}_n$ be the labels assigned to the n observations by a clustering algorithm. We denote by Σ the set of all possible permutations of the set of labels. The misclassification error rate, also called the “matching error” is defined as follows:

$$MCE = \min_{\sigma \in \Sigma} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{y_i \neq \sigma(\hat{y}_i)\}}$$

Recent approaches

DBSCAN

Density-based spatial clustering of applications with noise [4], a density-based, non-hierarchical algorithm.

- Fix two parameters: a reachability distance ϵ and a minimal number of points $MinPts$.
- Choose an arbitrary point P from the data and compute its ϵ -neighborhood, $V_\epsilon(P)$.
- If there are at least $MinPts$ in $V_\epsilon(P)$, point P (said to be dense) and all its neighbors will form a cluster C , otherwise it is labeled as noise.
- Any point $P' \in V_\epsilon(P)$ which is dense, is also added to the cluster C , together with its own ϵ -neighborhood.
- Once no dense points are found, a new unvisited point is retrieved and processed to explore new clusters.

Manhattan distance is often used as the dissimilarity measure in the ordinal case and the mutual information in the nominal case. See the *dbSCAN* function of the *fpc* package [2] from **R**.

Recent approaches

COBWEB ([5])

- A conceptual clustering approach based on the category utility (*CU*) measure.
- An incremental hierarchical algorithm is designed for qualitative data.
- Constructs a tree dynamically, inserting one individual at a time in the tree construction.
- At each individual insertion, four options are available options: inserting the individual to an existing cluster, creating a new cluster, merging two nodes, or splitting one node.
- Each option gives rise to a different partition whose *CU* is computed: the option maximizing the *CU* is selected.
- Loop over observations from the dataset, and stop when a minimal value of the category utility can not be exceeded.

The *Cobweb* function [1] of the *RWeka* package from R can be used.

Recent approaches

DIVCLUS-T

DIVCLUS-T [2] is a top-down hierarchical clustering method. It is monothetic, i.e subsets of observations in the dataset are split using single variable. Its goal is to optimize the same criterion as for classical CAH using Ward's method. It is designed for quantitative, qualitative and mixed data. To compute the splitting criterion, i.e the within-class inertia, DIVCLUS-T uses the euclidean distance in the quantitative case while it uses the Khi-2 distance in the qualitative case. The main advantage of this method, in comparison with other hierarchical methods, is that in addition to providing a dendrogram, each node is labeled by its binary splitting rule. So the dendrogram provided by DIVCLUS-T can be read exactly like a binary decision tree.

The algorithm is available with the *divclust* function of the *divclust* package from R.



Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS



Images

- 1 faire des classes d'images en utilisant toutes les méthodes de clusterisation non supervisé pour faire des classes
- 2 SEGMENTATION D'IMAGES: Couper l'image en zones homogènes.

Clustering de pixels à base de densité 1

Histogram-based methods (Wikipedia sur la segmentation d'image à base d'histogramme)

- Sont des méthodes très efficaces ne demandant qu'une seule passe au travers des pixels, pour calculer l'histogramme
- Ses pics et ses vallées servent à localiser les classes se basant soit sur la couleur soit sur l'intensité
- Peut être utilisé itérativement pour créer des classes plus fines
- Les pics et vallées peuvent être difficiles à déterminer
- Utilisables pour des multi-frames servant en vidéo pour suivre des mouvements

Dual Clustering

Méthode elle aussi basée sur l'histogramme des pixels (Wikipedia sur la segmentation d'image à base d'histogramme)

La méthode combine 3 caractéristiques de l'image La partition de l'image est construite sur les pics et les vallées de l'histogramme L'homogénéité de la partition est contrôlée par la compacité des clusters et la taille des gradients sur les bords des clusters.

Deux espaces sont introduits

- Le 1^{er} contient l'histogramme de brightness $H = H(B)$, qui permet de définir le seuil T de clusterisation en 2 classes (blanc,noir)
- le 2^{er} est l'espace dual de l'image originelle $B = B(x,y)$ à 3-dimensions. qui contient la bitmap-de-l'image est $b = \Phi(x,y)$, where $\Phi(x,y) = 0$, if $B(x,y) < T$, and $\Phi(x,y) = 1$, if $B(x,y) \geq T$.]

Une statistique mesure la compacité du bitmap.

Il reste à définir de bons bords en choisissant T .

Le seuil T étant donnée, on pose $M(T) = G/(k - L)$ has to be calculated where

- k est la différence entre la brightness de la classe et celle du background,
- L est la longueur du bord, G la moyenne du gradient sur le bord.on the borders.

$T^{opt} = \text{argmax}_T M(T)$ est le seuil optimal en 2 classes.

Annexe. Simu. and Perf. Eval.

LC simulation, Model M1

We set: $K = 3$, $p = 9$, $n \in \{100, 300, 500, 1000\}$, each variable has $m = 5$ levels, and clusters are equally sized.

Three clusters are defined, each characterized by a high frequency of one level: Level 1 (resp. 3 and 5) is the most frequent for cluster 1 (reps. 2 and 3). The other levels are uniformly distributed.

For example, The distribution for each variable in cluster one is given by:

$$\begin{aligned} P(X_j = 1) &= q \\ P(X_j = l) &= \frac{1-q}{m-1} \quad \text{for} \quad l \neq 1 \end{aligned}$$

We fix $q = 0.8$. This means, for each variable, the distribution of the five levels is: $(0.8, 0.2, 0.2, 0.2, 0.2)$ for the first cluster, $(0.2, 0.2, 0.8, 0.2, 0.2)$ for the second and $(0.2, 0.2, 0.2, 0.2, 0.8)$ for the third.

The same is done for the other clusters.

Annexe. Simu. and Perf. Eval.

A tree model, M3

$p = 3$, $k = 4$. Each variable X_j , $j \in \{1, \dots, p\}$, has $m = 6$ levels. Each level is coded as an integer, and we distinguish odd and even levels. The partition used for the simulation is shown in figure 9. Clusters are defined as follows:

- C1: x_1 and x_2 have odd levels, and x_3 is arbitrary
- C2: x_1 has odd levels, x_2 has even levels, and x_3 is arbitrary
- C3: x_1 has even levels, x_3 has odd levels, and x_2 is arbitrary
- C4: x_1 and x_3 have even levels, and x_2 is arbitrary

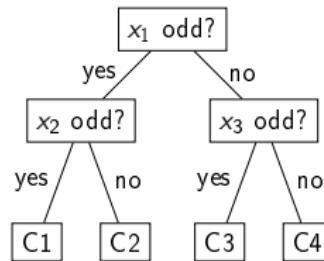


Figure: Tree structure used for data simulation model M3.

Annexe. Simu. and Perf. Eval.

Another tree model, M4

$p = 3$, $k = 4$. Here, each variable X_j , $j \in \{1, \dots, p\}$, has $m = 4$ levels. The only difference with M3 is that variable levels are not uniformly distributed in each cluster. Here, we consider a parameter p_0 that controls the non-uniformity of the distribution of levels, for example $p_0 = 0.8$. Clusters are defined as follows:

- C1: x_1 and x_2 have odd levels with $P(x_1 \in \{1, 3\}) = P(x_2 \in \{1, 3\}) = p_0$, and x_3 is arbitrary
- C2: x_1 has odd levels, x_2 has even levels with $P(x_1 \in \{1, 3\}) = P(x_2 \in \{2, 4\}) = p_0$, and x_3 is arbitrary
- C3: x_1 has even levels, x_3 has odd levels with $P(x_1 \in \{2, 4\}) = P(x_3 \in \{1, 3\}) = p_0$, and x_2 is arbitrary
- C4: x_1 and x_3 have even levels with $P(x_1 \in \{2, 4\}) = P(x_3 \in \{2, 4\}) = p_0$, and x_2 is arbitrary



Annexe. Simu. and Perf. Eval.

Les méthodes de cette partie

- ACP
- ICA comme extention
- Multidimensional Scaling de tableau de distances

| | | | | | | | |
|---------|-------------------|----------------------------|---|------------|----------------|------------|-----|
| Densité | Clustering ooo | Clustering. Images oooo | ACP et la famille  | ICA ooo | MDS ooooooo | Ressembler | TP3 |
|---------|-------------------|----------------------------|---|------------|----------------|------------|-----|

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS



ACP et la famille

Quelle est la question? ACP/ICA (1)

Des mesures (p) sont faites à l'aide de descripteurs sur des individus (n). Certaines sont utiles, d'autres superflues.

- ACP: construire un petit nombre de **variables NOUVELLES** qui discriminent au mieux les descriptions des individus. Elles seront appelées les Composantes principales
- ICA comme extension: **séparer les sources.**

ACP et la famille

Quelle est la question? MDS (2)

- Multidimensional Scaling de tableau de dis-similarités. La connaissance de départ est un indice de dis-similarité entre les couples de descripteurs initiaux, en général des distances (plus la distance est grande plus les différences sont grandes). Appelons DIS ces informations.
 Le problème est de construire de nouveaux points (appelés $(i^{repr}, i = 1 \dots n)$) admettant une description simple couplée aux individus ($i, i = 1 \dots n$), et pour lesquels il sera facile de lire sur des figures la dissimilarité DIS^{repr} . On veut que ces nouveaux points possèdent la propriété $DIS^{repr} - DIS^{repr}$ soit **globalement petits**.

ACP et la famille

PCA.Hastie 1

PCA produces a low-dimensional representation of a dataset.
It finds a sequence of linear combinations of the variables that

- 1 have maximal variance**
- 2 mutually uncorrelated.**

PCA also serves as a tool for data visualization and understanding

The 1rst principal component of a set of features

X_1, X_2, \dots, X_p is the normalized linear combination of the features

$$Z_1 = \phi_{11}X_1 + \phi_{21}X_2 + \cdots + \phi_{p1}X_p$$

that has the largest variance. Moreover $\sum_{j=1}^p (\phi_{j1})^2 = 1$

The $\phi_1 = (\phi_{11}, \phi_{21}, \dots, \phi_{p1})'$ are the **loadings** defining the **1rst principal component loading vector**,

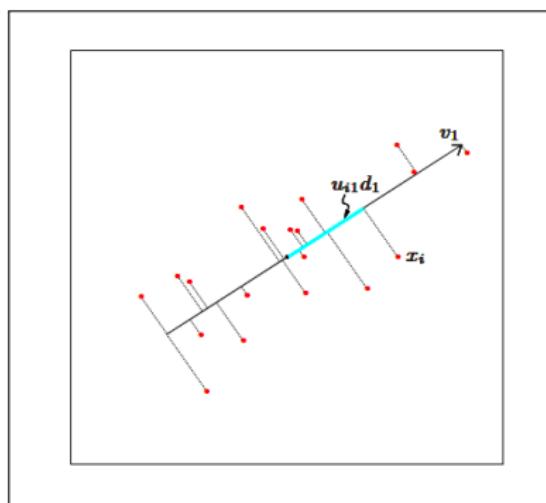
ACP et la famille

Mettre un tableau $X_{n,p}$; Nuage

- 1 $X_{n,p}, n = 8, p = 2$
- 2 Tableau centré $scale(X)$ moyenne des colonnes = 0, var des colonnes = 1.
- 3 Tableau de correlation $R_{p,p} = (1/n)X'X$;
- 4 $R = \begin{pmatrix} 1 & r \\ r & 0 \end{pmatrix}$
- 5 $V = (1/\sqrt{2}) \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$
- 6 $RV_1 = (1+r)V_1; RV_2 = (1-r)V_1;$
- 7 $CompoPrincip_1 = XV_1, CompoPrincip_2 = XV_2$ sont de moyenne nulle
- 8 $var(CompoPrincip_1) = (1/n)XV_1, CompoPrincip_2 = XV_2$

ACP et la famille**Nuage projeté 1D**
 $a_i = \text{Projection}_{[a]}(X_i)$ Variance et Inertie

$\text{var}(x) = (1/n) \sum_{i=1}^n (x_i - \bar{x})^2$



Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

○○○

○○○○

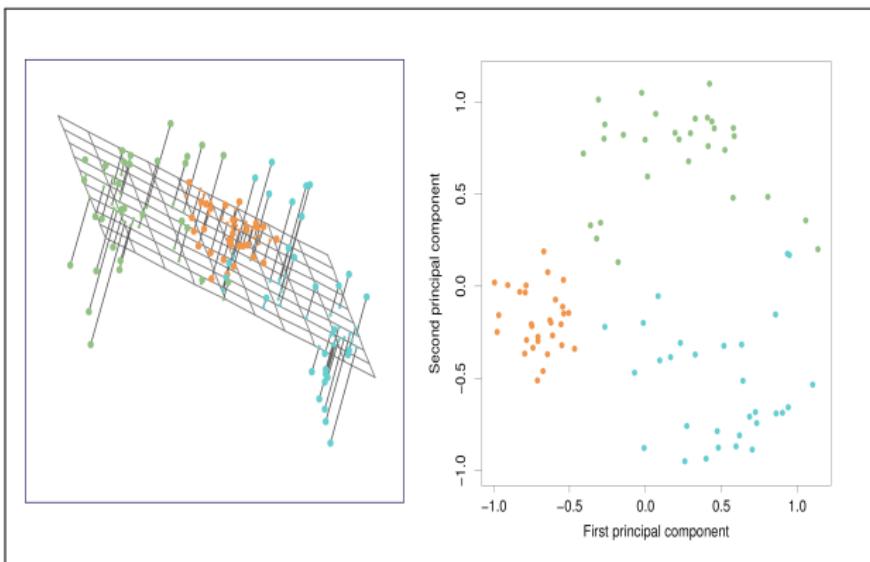


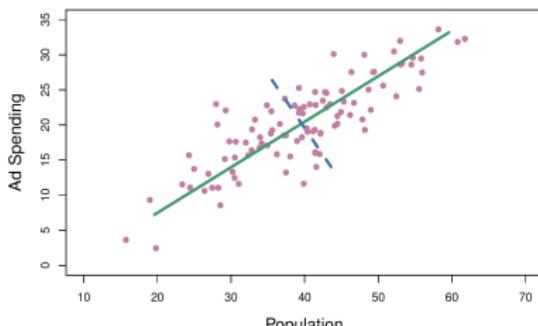
○○○

○○○○○○○

ACP et la famille

Nuage projeté 2D



ACP et la famille**PCA Un axe. Intuition**

The population size (`pop`) and ad spending (`ad`) for 100 different cities are shown as purple circles. The green solid line indicates the first principal component direction, and the blue dashed line indicates the second principal component direction.

Figure: HastiePCA7



Exemple: Notes

8 élèves, 4 matières de dimension 2 (1)

On a enregistré les notes de 8 élèves dans 4 matières mathématique (V1), physique (V2), Français (V3), latin (V4).

Une ACP centrée réduite recherche les associations entre les matières.

- Le premier axe est celui du **facteur général**, regroupement des 4 notes de chacun des élèves. Les coefficients sont ceux de la première colonne du tableau des vecteurs propres. Sur ce premier axe, l'ordre des élèves est celui des moyennes.
- Le second axe oppose **les matières scientifiques V1 et V2 à V3 et V4, variables littéraires**. Les élèves ayant de bons résultats dans les matières littéraires s'opposent à ceux dont les résultats sont bons dans les matières scientifiques.
- La **dimension** du problème est égale à 2.

Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

ooo

oooo



ooo

oooooo

Exemple: Notes

(2)



Données 8 élèves et des notes

| | | | | |
|---|----|----|----|----|
| 1 | 13 | 12 | 08 | 09 |
| 2 | 14 | 14 | 15 | 15 |
| 3 | 05 | 07 | 14 | 11 |
| 4 | 14 | 14 | 12 | 12 |
| 5 | 11 | 10 | 05 | 07 |
| 6 | 08 | 08 | 08 | 08 |
| 7 | 06 | 07 | 11 | 09 |
| 8 | 06 | 06 | 05 | 05 |

ooo

oooo

```
ooooooo
oo●ooooo
oooooooo
oooo
oooooo
```

ooo

ooooooo

Exemple: Notes

(3)

 Résultats

| | |
|--|---|
| matrice à diagonaliser | vecteurs propres en colonne |
| <pre>===== 1.0000 0.9841 0.2503 0.5430 0.9841 1.0000 0.3980 0.6655 0.2503 0.3980 1.0000 0.9484 0.5430 0.6655 0.9484 1.0000 =====</pre> | <pre>===== 0.48 -0.54 0.62 -0.28 0.52 -0.41 -0.74 0.05 0.44 0.63 -0.08 -0.6 0.54 0.37 0.23 0.71 =====</pre> |
| valeurs propres | composantes principales |
| <pre>===== 2.9119 1.0836 0.0040 0.0005 =====</pre> | <pre>===== 0.6657 -1.1456 -0.0682 -0.0367 3.0155 0.3667 0.1144 0.0013 -0.2424 1.9534 -0.0271 -0.0070 2.0498 -0.4955 -0.0932 0.0376 -0.8704 -1.3581 0.0505 -0.0126 -1.1250 -0.0667 0.0280 -0.0005 -0.9178 1.0918 -0.0309 -0.0147 -2.5754 -0.3460 0.0265 0.0324 =====</pre> |
| qualité de l ajustement | |
| <pre>===== 72.79 99.88 99.98 100.00 ===== ==</pre> | |

Ondelettes et applications Debruitage 3

Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

○○○

○○○○

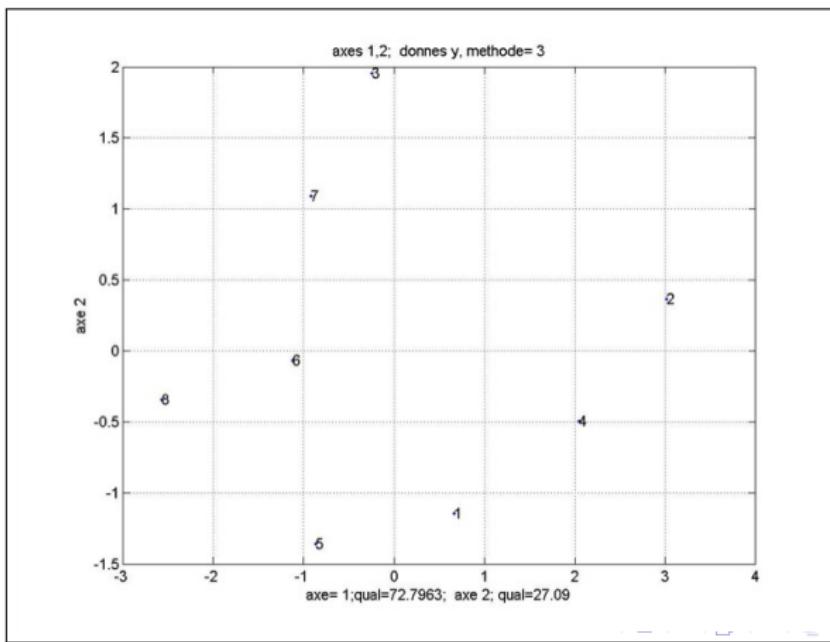


○○○

○○○○○○

Exemple: Notes

(4)



Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

○○○

○○○○

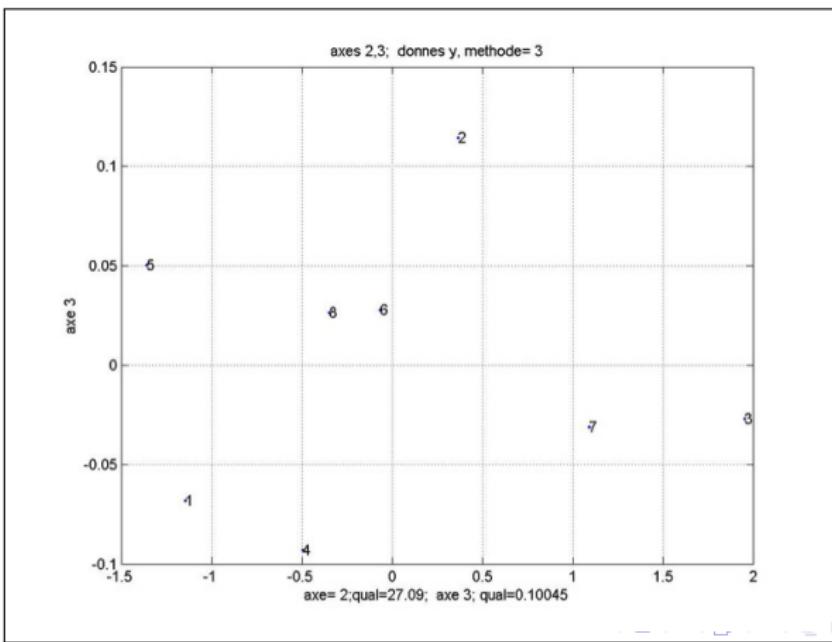


○○○

○○○○○○

Exemple: Notes

(5)



Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

```

    oooooo
    ooooo●o
    oooooo
    oooo
    oooooo
  
```

ICA

○○○

MDS

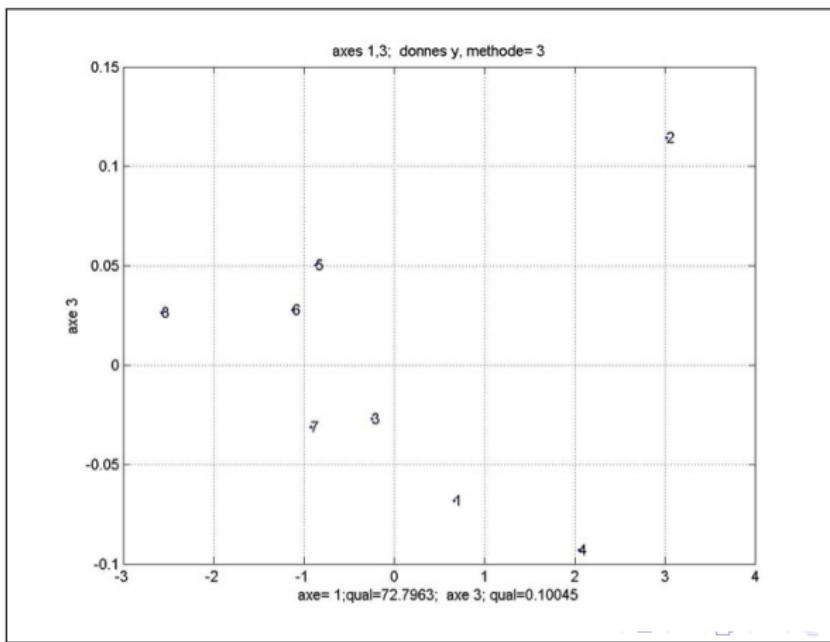
○○○○○○

Ressembler

TP3

Exemple: Notes

(6)



Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

```

○○○○○○○
○○○○○●○○
○○○○○○○
○○○○○
○○○○○

```

ICA

○○○

MDS

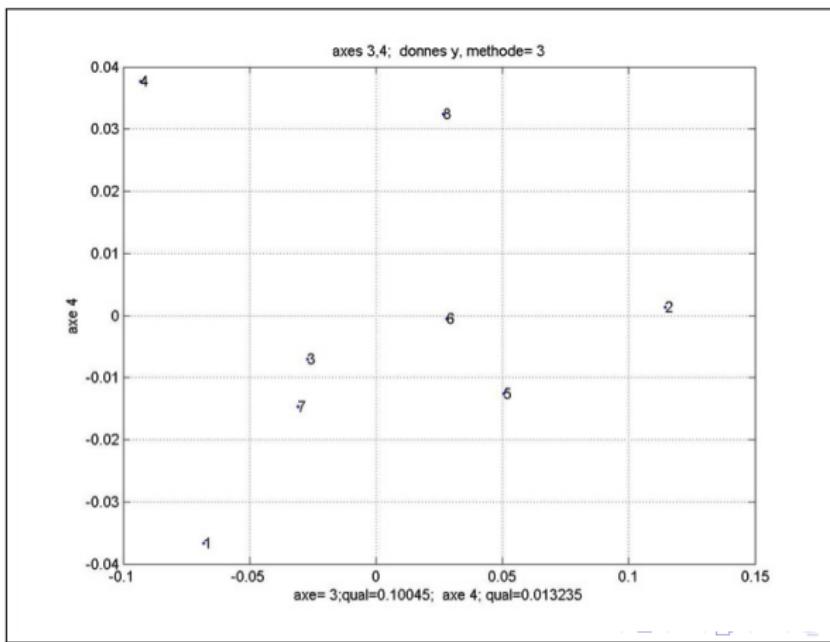
○○○○○○○

Ressembler

TP3

Exemple: Notes

(7)



Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

○○○

○○○○

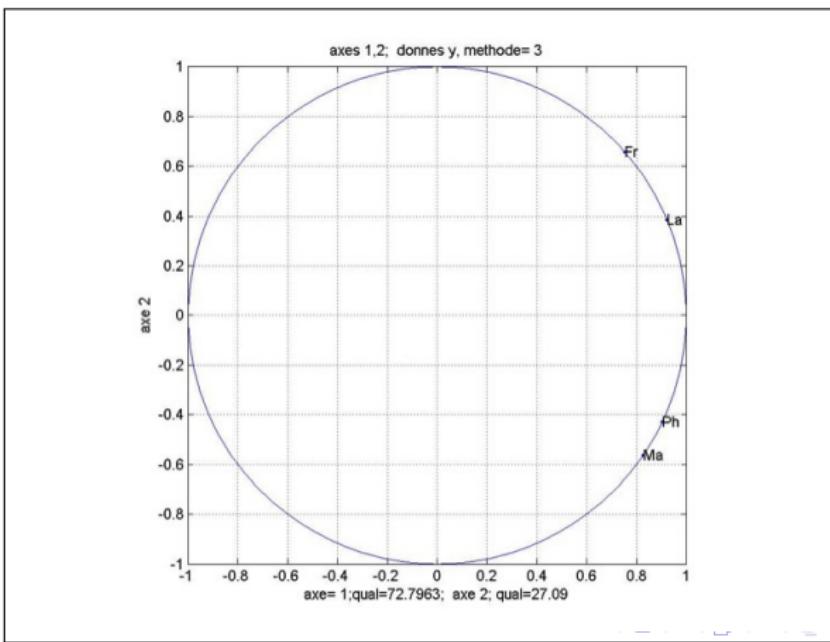


○○○

○○○○○○

Exemple: Notes

(8)



Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

○○○

○○○○

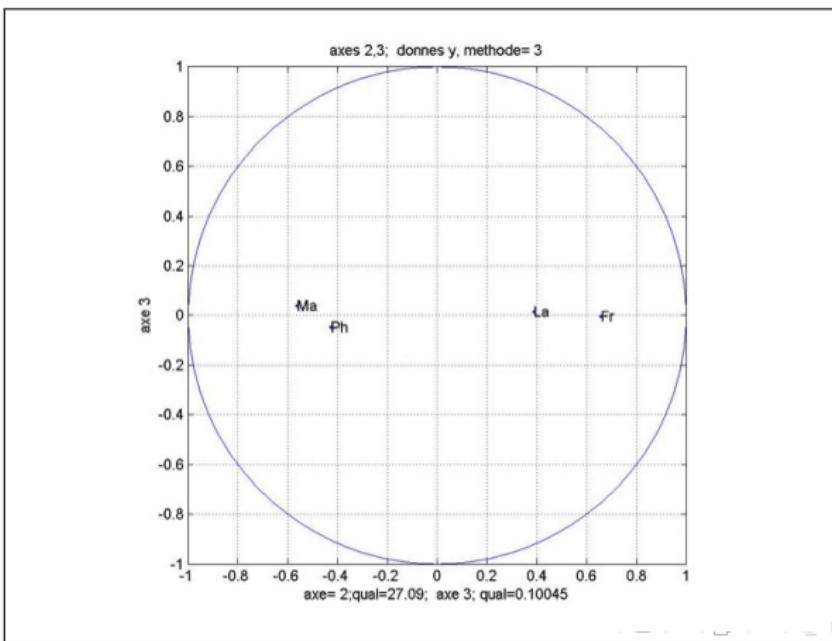


○○○

○○○○○○

Exemple: Notes

(9)



**ACP chiens et loups**

Les données des chiens

Les ossements de cranes de 42 canidés, identifiés comme chiens et loups, ont été mesurés. Les variables sont : Les lignes correspondent aux **individus** et les colonnes aux **variables**

Le nombre d'individu est noté *n* et le nombre de colonnes est noté *p*

- LCB : longueur condylo-basale
- LSM : longueur de la mâchoire supérieure
- LBM : largeur bi-maxillaire
- LP : longueur de la carnassière supérieure
- LM : longueur 1^{er} molaire supérieure
- LAM : largeur 1^{er} molaire supérieure

| TYPE | LCB | LSM | LBM | LP | LM | LAM |
|------------------|------|------|-----|-----|-----|-----|
| BULL-DOG 1 | 1290 | 640 | 950 | 175 | 112 | 138 |
| BULL-DOG 2 | 1540 | 740 | 760 | 200 | 142 | 165 |
| BOXER | 1580 | 710 | 710 | 167 | 125 | 133 |
| SAIN T-BERNARD | 2200 | 1110 | 880 | 225 | 154 | 180 |
| BULL-MASSIF | 1900 | 930 | 780 | 197 | 132 | 140 |
| DOGUE ALLEMAND 1 | 2410 | 1190 | 870 | 210 | 147 | 183 |
| DOGUE ALLEMAND 2 | 2420 | 1200 | 850 | 199 | 153 | 176 |
| SETTER 1 | 2010 | 1050 | 700 | 198 | 143 | 159 |
| SETTER 2 | 1960 | 1060 | 670 | 185 | 126 | 142 |

Repris de Cours d'ACP. 2010. Nicolas Durrande - durrande@emse.fr. Script

en R Cours 1 ACP. Repris de Durande 2010. Source : M.JAMBU - Classification

Densité

Clustering

Clustering. Images

ACP et la famille

ICA

MDS

Ressembler

TP3

○○○

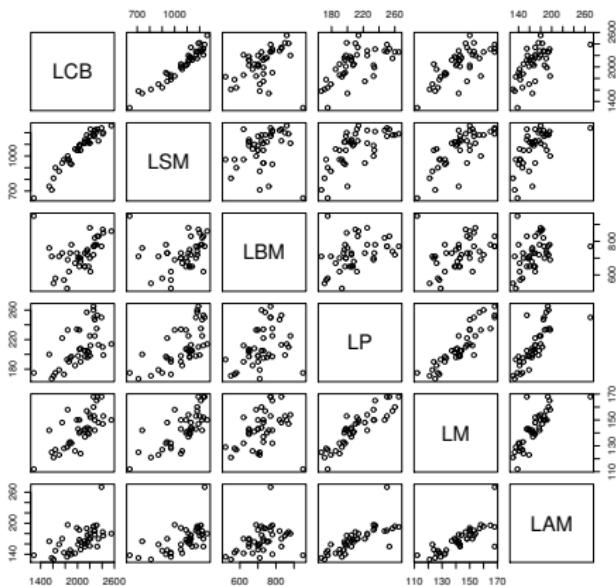
○○○○

○○○

○○○○○○

ACP chiens et loups

Exemple sur les données des chiens



ACP chiens et loups

Distances et Inertie

Dissance et le produit scalaire euclidiens et inneriens sont définis par :

$$\begin{aligned} d^2(x_i, x_j) &= ||x_i - x_j||^2 = \sum_{k=1}^p (x_i^k - x_j^k)^2 \\ \langle x_i, x_j \rangle &= \sum_{k=1}^p x_i^k x_j^k = X_i' X_j \end{aligned}$$

Inertie par rapport à un point a: $I_a = \sum_{i=1}^n \frac{1}{n} d^2(x_i, a)$

Or $I_a = I_g + ||a - g||^2$. L'inertie du siège par rapport à un set $E \subset R^p$:

$I_E = \sum_i \frac{1}{n} d^2(x_i, E) = \sum_i \frac{1}{n} d^2(x_i, P_E(x_i))$

Soit C la matrice de covariance empêtrée des variables

$$\begin{aligned} C &= \sum_{i=1}^n \frac{1}{n} (x_i - g)^T (x_i - g) \\ C &= P D P^T \\ f_1 &= e_1 P, ||f_1|| = 1 \\ I_{\Delta_1^{\perp}} &= e_1 P D P^T e_1^T = f_1 D f_1^T = \sum_{i=1}^p \lambda^i (f_i^1)^2 \\ &\leq \max(\lambda^i) \sum_{i=1}^p (f_i^1)^2 = \max(\lambda^i) \end{aligned}$$

Le maximum de $I_{\Delta_1^{\perp}}$ est atteint pour $f_1 = (1, 0, \dots, 0)$, $e_1 = f_1 P^T$ donc e_1 est le premier vecteur propre de C.

Le premier VP de C est appellé le 1^{er} axe factoriel ou la 1^{ere} direction principale.

Les valeurs des projections des données suivant ces axes s'appellent les composantes principales.

On cherche ensuite un axe Δ_2 posé par e_2 , orthogonal à Δ_1 tel que l'inertie d'a siège par rapport à plus (g, e_1, e_2) soit minimale.

$$\begin{aligned} \Delta_2 &= \operatorname{argmax}(I_{(g \oplus e_2)^{\perp}}) \\ &= \operatorname{argmax}(\sum_i \frac{1}{n} d^2(x_i, (\Delta_1 \oplus \Delta_2)^{\perp})) \\ &= \operatorname{argmax}(\sum_i \frac{1}{n} d^2(x_i, g) - d^2(x_i, \Delta_1) - d^2(x_i, \Delta_2)) \end{aligned}$$

Étiquetage :

Les k premières directions principales correspondent aux VP de la matrice de covariance associées aux k plus grandes vp. Elles sont deux à deux orthogonales.

Les plus des 2 axes factoriels sont les plus factoriels. Les variables associées aux axes sont le CP.



ACP chiens et loups

Retour sur l'exemple des chiens-loups

Peut-on interpréter ces résultats ?

- Quels sont les liens entre les nouvelles variables et les anciennes ?
- Les nouvelles variables ont elle un sens ?

On va pour cela s'intéresser à la corrélation entre composantes principales et variables.



ACP chiens et loups

Retour sur l'exemple des chiens-loups

Notons $X_c = X - g$ la matrice des données centrées. Les composantes principales Y^1, Y^2, \dots, Y^P sont donc les colonnes de la matrice $Y = X_c P$. On en déduit que la corrélation entre X^i et Y^j vaut

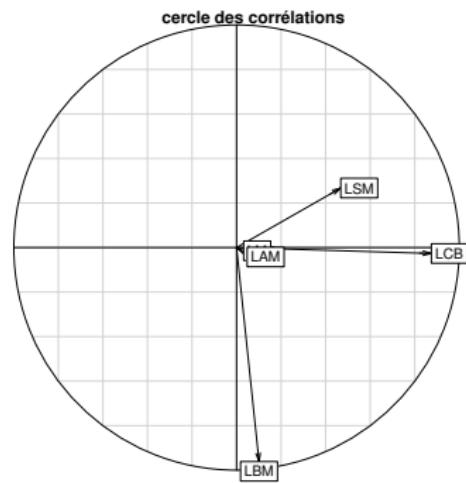
$$\text{corr}(X^i, Y^j) = \frac{\sqrt{\lambda_j}}{\sqrt{\text{var}(X^i)}} p_i^j$$

On peut alors représenter les résultats que l'on vient d'obtenir sur un graphique appelé **cercle des corrélations**.

**ACP chiens et loups**

Retour sur l'exemple des chiens-loups

Cercle des corrélations





ACP chiens et loups

Problèmes posés

- Les résultats dépendent de la variance des variables de départ
- Si les variables non homogènes, il y a un problème de sens physique

Pour s'affranchir de ces deux problèmes, on va utiliser la méthode d'ACP normée.

ACP normée

Centrer et réduire une variable=soustrait sa moyenne et on la divise par son écart type : $X_{cr} = \frac{X - mean(X)}{sd[X]}$. Cela répond aux deux problèmes que l'on avait.

- Les résultats ne dépendent pas de la variance des variables
- Les données sont adimensionnées

Centrer-réduire les variables correspond à un changement de métrique.

$$\tilde{x} = x - mean(x)$$

$$d^2(x_i, x_j) = \sum_{k=1}^p \frac{(\tilde{x}_i^k - \tilde{x}_j^k)^2}{var(x^k)}. \quad (1)$$

Si on reprend la démonstration de l'ACP, les axes principaux sont les vecteurs propres de la **matrice R de corrélation**.

- L'inertie totale est égale à p
- Le cercle des corrélations représente les variables dans la base des directions propres

l'ACP centrée-réduite assure que les variables présentant de faibles variances ne sont pas "écrasées" par les autres.

Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

ICA

○○○

MDS

○○○○○○○

Ressembler

TP3

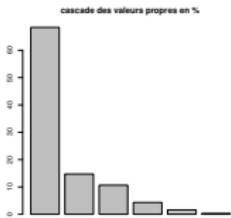
ACP normée

Retour sur l'exemple des chiens-loups

$$\text{Cov} = \begin{pmatrix} 76467 & 40311 & 8899 & 4463 & 2718 & 4050 \\ 40311 & 23020 & 2803 & 2632 & 1524 & 2249 \\ 8899 & 2803 & 8522 & 896 & 441 & 816 \\ 4463 & 2632 & 896 & 689 & 320 & 499 \\ 2718 & 1524 & 441 & 320 & 186 & 268 \\ 4050 & 2249 & 816 & 499 & 268 & 621 \end{pmatrix}$$

La matrice de corrélation :

$$\text{Corr} = \begin{pmatrix} 1.00 & 0.96 & 0.34 & 0.61 & 0.71 & 0.58 \\ 0.96 & 1.00 & 0.20 & 0.66 & 0.73 & 0.59 \\ 0.34 & 0.20 & 1.00 & 0.36 & 0.35 & 0.35 \\ 0.61 & 0.66 & 0.36 & 1.00 & 0.89 & 0.76 \\ 0.71 & 0.73 & 0.35 & 0.89 & 1.00 & 0.78 \\ 0.58 & 0.59 & 0.35 & 0.76 & 0.78 & 1.00 \end{pmatrix}$$



Densité

Clustering

○○○

Clustering. Images

○○○○

ACP et la famille

ICA

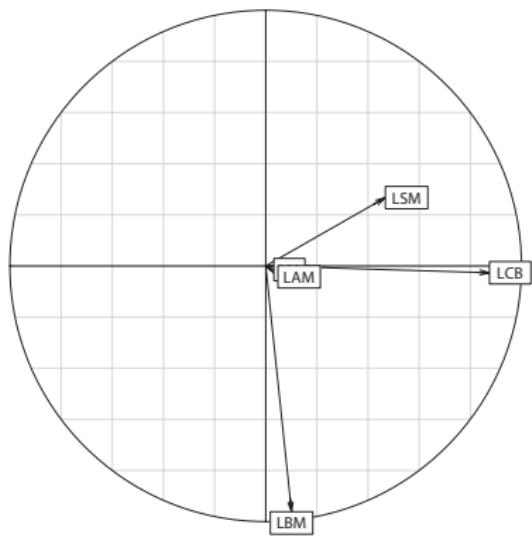
○○○

MDS

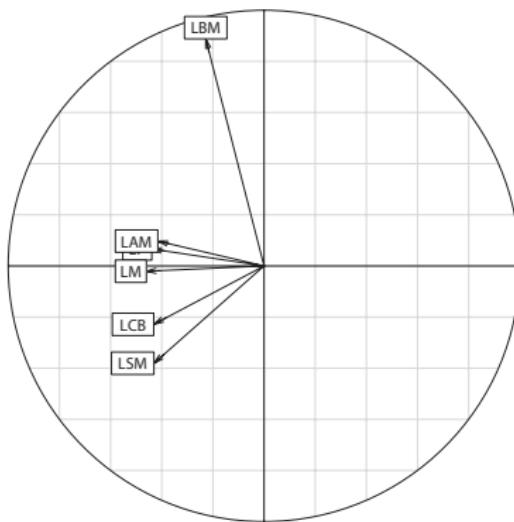
○○○○○○

Ressembler

TP3

ACP normée**Directions trouvées**

Cercle des corrélations ACP classique

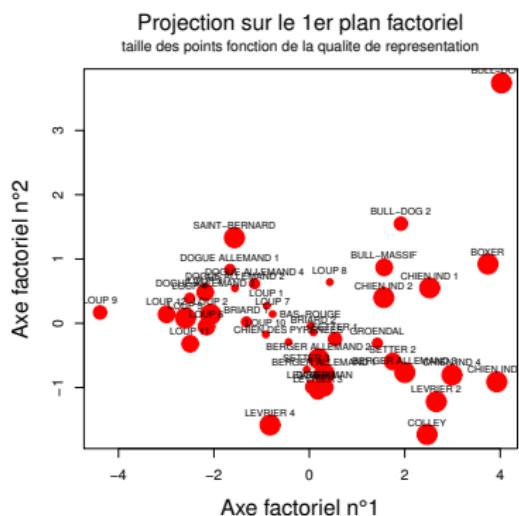


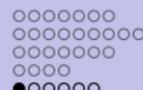
Cercle des corrélations ACP normée

ACP normée

Retour sur l'exemple des chiens-loups

Dans l'espace des individus





Exemple 2

Conditions de travail et au niveau de vie dans différentes villes du monde.

On s'intéresse aux conditions de travail et au niveau de vie dans différentes villes du monde. On étudie $n = 46$ villes sur lesquelles on mesure $p = 3$ variables quantitatives :

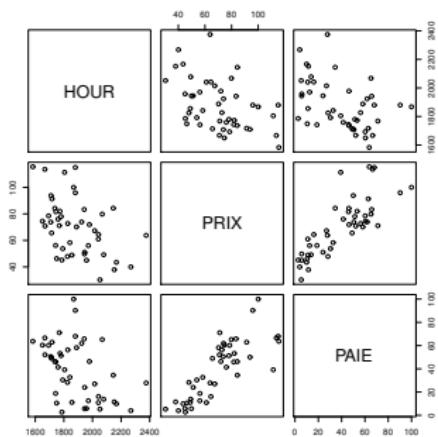
- HOUR : nb moyen d'heures de travail dans 12 activités
- PRIX : indice des prix sur la base de 112 produits et services
- PAIE : index des salaires horaires dans 12 activités, déductions faites

source : *Prices and earnings around the globe Economic Research Department, Union Bank of Switzerland, Zurich.*

| Ville | HOUR | PRIX | PAIE |
|--------------|------|-------|-------|
| Amsterdam | 1714 | 65.6 | 49.0 |
| Athens | 1792 | 53.8 | 30.4 |
| Bogota | 2152 | 37.9 | 11.5 |
| Bombay | 2052 | 30.3 | 5.3 |
| Brussels | 1708 | 73.8 | 50.5 |
| Buenos Aires | 1971 | 56.1 | 12.5 |
| Caracas | 2041 | 61.0 | 10.9 |
| Chicago | 1924 | 73.9 | 61.9 |
| Copenhagen | 1717 | 91.3 | 62.9 |
| Dublin | 1759 | 76.0 | 41.4 |
| Dusseldorf | 1693 | 78.5 | 60.2 |
| Frankfurt | 1650 | 74.5 | 60.4 |
| Geneva | 1880 | 95.9 | 90.3 |
| ... | ... | ... | ... |
| Zurich | 1868 | 100.0 | 100.0 |

Exemple 2

Diagramme pairs et matrice de corrélation



Matrice de corrélation :

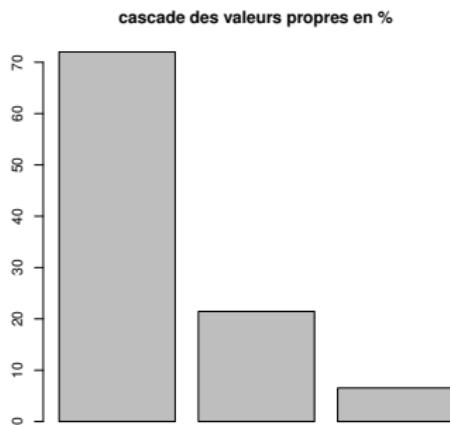
| | HOUR | PRIX | PAIE |
|------|-------|-------|-------|
| HOUR | 1.00 | -0.45 | -0.45 |
| PRIX | -0.45 | 1.00 | 0.80 |
| PAIE | -0.45 | 0.80 | 1.00 |



Exemple 2

VP et vp

Valeurs propres et vecteurs propres

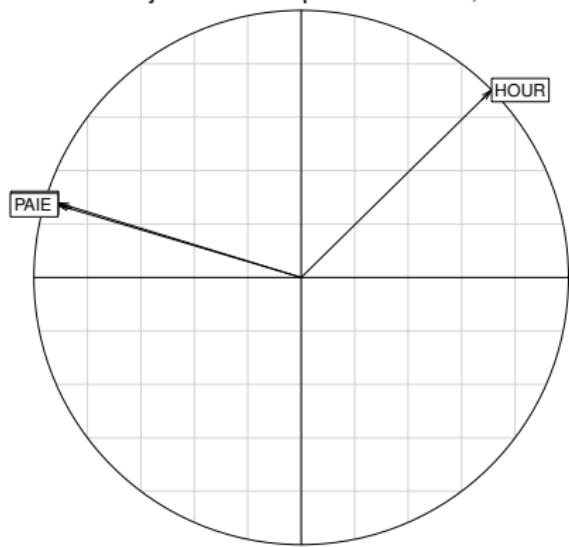


| | VP1 | VP2 | VP3 |
|------|-------|------|-------|
| HOUR | 0.48 | 0.87 | 0.00 |
| PRIX | -0.61 | 0.34 | -0.70 |
| PAIE | -0.61 | 0.33 | 0.70 |

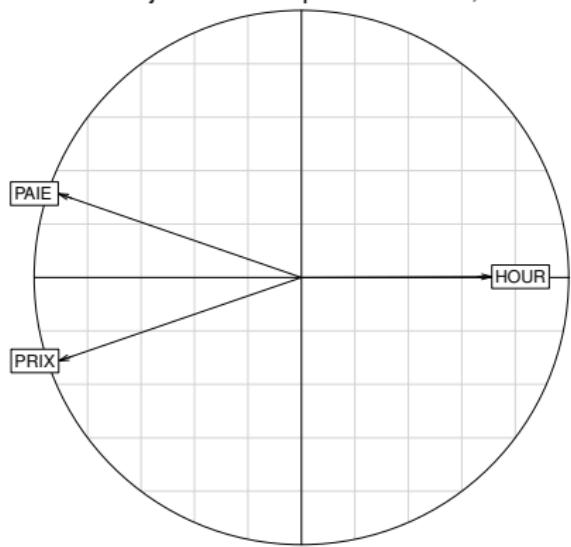


Exemple 2

Projection sur le plan factoriel 1,2



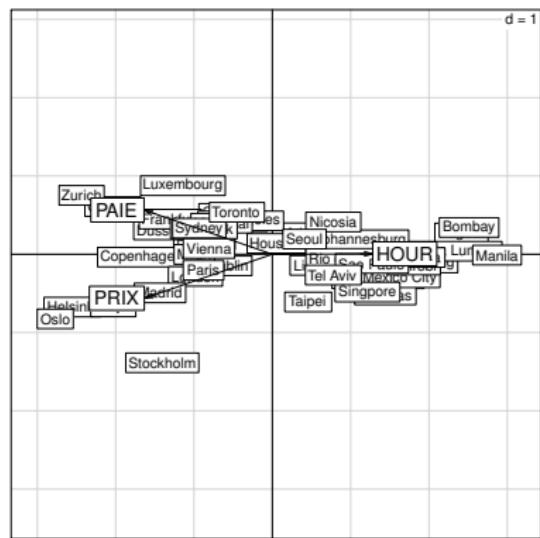
Projection sur le plan factoriel 1,3



Peut-on interpréter les directions que l'on a trouvées ?



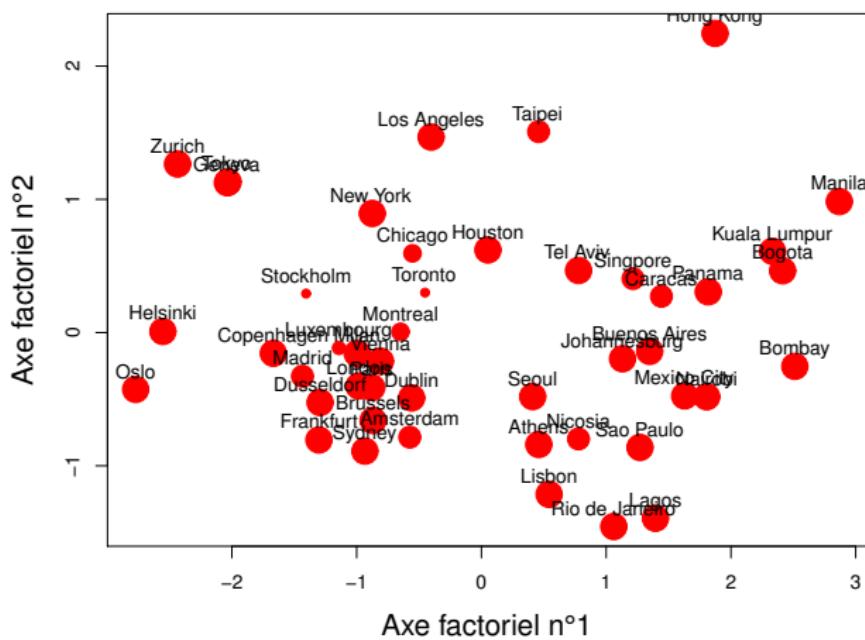
Exemple 2





Exemple 2

Projection sur le 1er plan factoriel
taille des points fonction de la qualité de représentation



| | | | | | | | |
|---------|-------------------|----------------------------|---|------------|----------------|------------|-----|
| Densité | Clustering ooo | Clustering. Images oooo | ACP et la famille  | ICA ooo | MDS ooooooo | Ressembler | TP3 |
|---------|-------------------|----------------------------|---|------------|----------------|------------|-----|

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

ICA: Analyse en Composantes Indépendantes

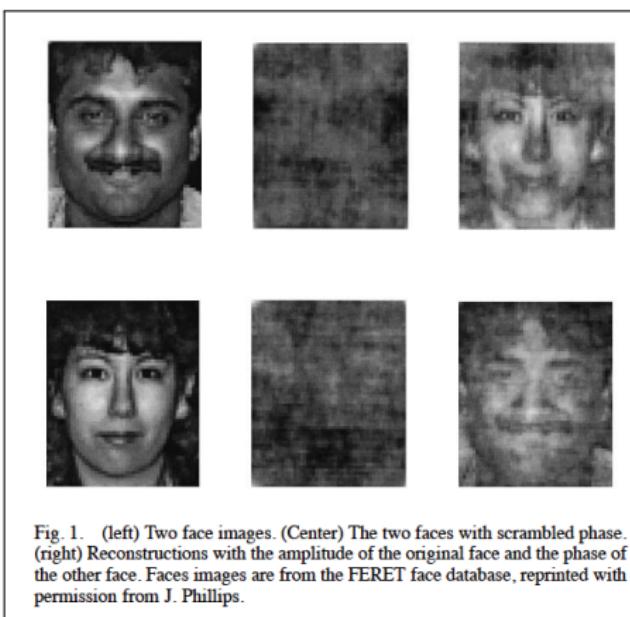
- voir [*https://fr.wikipedia.org/wiki/Analyse_en_composantes_independantes*](https://fr.wikipedia.org/wiki/Analyse_en_composantes_independantes)
- Etudie le même problème que l'ACP mais remplace l'obligation de n'être pas corrélée par l'obligation d'être indépendant qui est plus forte.
- ACP sans sous-basement gaussien ni quadratique
- Basé sur l'entropie mutuelle
- L'illustration classique de la séparation de sources est le problème suivant:
Lors d'une soirée, on dispose de P microphones dans une salle où N personnes discutent par groupes de tailles diverses.
 - Chaque microphone enregistre la superposition des discours des personnes à ses alentours et
 - le problème consiste à retrouver la voix de chaque personne débarrassée des autres voix considérées comme parasites.
- $x = As + \epsilon$ où les inconnues sont les sources s .



ICA: Analyse en Composantes Indépendantes

- [Cardoso J.F] Il s'agit de méthodes traitant des observations multivariées) afin d'en extraire des composantes linéaires qui soient aussi indépendantes que possible. Cette simple idée s'est révélée très fructueuse pour le traitement des signaux multi-capteurs dans de nombreux domaines : réseaux d'antennes pour les télécommunications, prise de son, signaux biomédicaux multi électrodes, et de manière plus générale dans tous les cas où un système de plusieurs capteurs, fournissant des signaux cohérents, est à l'écoute d'un ensemble discret de sources de signal que l'on cherche à extraire des observations et que l'on peut, pour des raisons physiques, supposer mutuellement statistiquement indépendantes.

ICA: Analyse en Composantes Indépendantes.Image



| | | | | | | | |
|---------|-------------------|----------------------------|--|------------|-----------------|------------|-----|
| Densité | Clustering ooo | Clustering. Images oooo | ACP et la famille oooooooo oooooooooooo oooooooooooo ooooo oooooo | ICA ooo | MDS oooooooo | Ressembler | TP3 |
|---------|-------------------|----------------------------|--|------------|-----------------|------------|-----|

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

MDS 1



But de MDS

On se donne des indices $\delta_{ii'}, 1 \leq i, i' \leq N$

On veut représenter N points,
dans un espace de petite
dimension

de sorte que les distances
**euclidiennes de la
représentation** soient aussi
proches que possible des
indices

$$d_{ii'}, 1 \leq i, i' \leq N$$

$$d_{ii'} \approx \delta_{ii'}$$

49



Critères pour le calcul. CRIT1

Trouver dans un espace de petite taille $z_i \in R^k$

Least squares ou Kruskall-Shepard scaling

STRESS function $S(z_1, \dots, z_n) = \left[\sum_{1 \leq i < j \leq n} (d_{ij} - \|z_i - z_j\|)^2 \right]^{1/2}$ Soit minimum

$$\text{ArgMin}_{z_1, \dots, z_n} S(z_1, \dots, z_n)$$

Par un algorithme de descente du gradient

50

MDS 2

Critère sur produits scalaires CRIT2



Nécessité des produits scalaires et pas seulement des distances

Classical scaling. Les données sont les points x_i

similarité $s_{ij} = \langle x_i - \bar{x}, x_j - \bar{x} \rangle$

$$\Rightarrow S(z_1, \dots, z_n) = \left[\sum_{1 \leq i < j \leq n} (s_{ij} - \langle z_i - \bar{z}, z_j - \bar{z} \rangle)^2 \right]^{1/2}$$

Solution: soit la diagonalisation $\underline{S} = V' D V$, z_1, \dots, z_n sont les lignes de

$$Z = V D V^{-1}$$

$$V_{i,k} \quad D_{i,i} \quad z_{i,k}$$

51

isoMDS en R

- isoMDS(MASS) **Kruskal's Non-metric Multidimensional Scaling**
- One form of non-metric multidimensional scaling

- isoMDS(d, y = cmdscale(d, k), k = 2, maxit = 50, trace = TRUE, tol = 1e-3, p = 2)

Arguments

- d distance structure of the form returned by dist, or a full, symmetric matrix. Data are assumed to be dissimilarities or relative distances, but must be positive except for self-distance. Both missing and infinite values are allowed.
- y An initial configuration. If none is supplied, cmdscale is used to provide the classical solution, unless there are missing or infinite dissimilarities.
- k The desired dimension for the solution, passed to cmdscale.
- maxit The maximum number of iterations.
- trace Logical for tracing optimization. Default TRUE.
- tol Convergence tolerance.
- p Power for Minkowski distance in the configuration space.

52

MDS 3



Critères pour le calcul CRIT3

Sammon mapping (conserve les petites valeurs)

$$\rightarrow S(z_1, \dots, z_n) = \left[\sum_{i \neq j} \frac{(d_0 - \|z_i - z_j\|)^2}{d_0} \right]^{1/2}$$

Favorise les petites distances
ce qui permet la représentation
des propriétés locales mieux que dans
les méthodes concurrentes

Exemple de R (complété)

Swiss Fertility and Socioeconomic Indicators (1888)

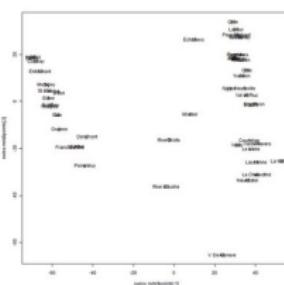
Standardized fertility measure and socio-economic indicators for each of
47 French-speaking provinces of Switzerland at about 1888.

A data frame with 47 observations on 6 variables, each of which is in %

- [,1] Fertility lg. 'common standardized fertility measure'
- [,2] Agriculture % of males involved in agriculture as occupation
- [,3] Examination % draftees receiving highest mark on army examination
- [,4] Education % education beyond primary school for draftees.
- [,5] Catholic % 'catholic' (as opposed to 'protestant').
- [,6] InfantMortality live birth who live less than 1 year.

53

Essai de isoMDS avec une distance euclidienne entre les points
de l'espace initial



```

data(swiss)
swiss.x <- acr.mds(swiss[, -1])
swiss.samp <- acr.mds(swiss)
swiss.dist <- isoMDS(swiss.dist, dist.k=3, p=2)
plot(swiss.mdspoints, type = "n");
text(swiss.mdspoints, labels = norm);
text(swiss.mdspoints, labels = norm)

swiss.ah <- Shepard(swiss.dat, swiss.mds(points))
plot(swiss.ah, pch = "+");
lines(swiss.ahx, swiss.mdsy, type = "S")

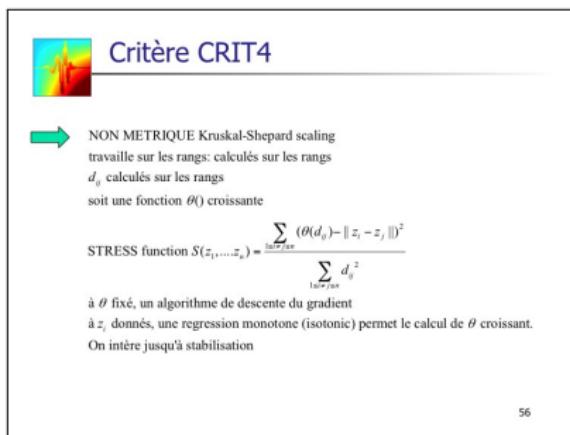
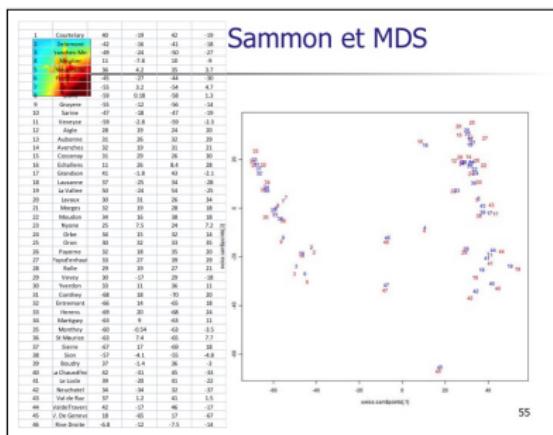
# -----
# sammon
plot(swiss.sampoints, type = "n");
text(swiss.sampoints, labels = norm);
text(swiss.sampoints, labels = norm);
text(swiss.sampoints, labels = norm);

# -----
# comparaison des deux analyses
plot(swiss.sampoints, type = "n");
text(swiss.sampoints, labels = norm);
text(swiss.sampoints, labels = norm);
text(swiss.sampoints, labels = norm);
text(swiss.mdspoints, labels =
  as.character(1:nrow(swiss.x)), col = blue)
text(swiss.mdspoints, labels =
  as.character(1:nrow(swiss.x)), col = blue)

```

54

MDS 4



MDS 5



Exemple Voitures mds

57



```
% Demos stats en MATLAB mdscale
load carbigr
X = [MPG, Acceleration, Displacement, Weight, Horsepower];
models77 = find((Model_Year==77));

h = glyphplot(X(1:9,:), 'glyph', 'star', 'varLabels', varNames, 'obslabels',
Model(1:9,:));

dissimilarity = pdist(zscore(X(models77,:)));
Y = mdscale(dissimilarity,2)
```

58

| Densité | Clustering ooo | Clustering. Images oooo | ACP et la famille oooooooo oooooooooooo oooooooooooo ooooo oooooo | ICA ooo | MDS oooooooo | Ressembler | TP3 |
|---------|-------------------|----------------------------|--|------------|-----------------|------------|-----|
|---------|-------------------|----------------------------|--|------------|-----------------|------------|-----|

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

| Densité | Clustering ○○○ | Clustering. ○○○○ | Images | ACP et la famille ○○○○○○ ○○○○○○○○ ○○○○○○ ○○○○○ | ICA ○○○ | MDS ○○○○○○ | Ressembler | TP3 |
|---------|-------------------|---------------------|--------|--|------------|---------------|------------|-----|
|---------|-------------------|---------------------|--------|--|------------|---------------|------------|-----|

Similarité ou dissimilarité

Objets entre lesquels une similarité est calculée

- Vecteur, Matrice, Cube, Courbe, Fonction, Surface, Image, Espace vectoriel, Opérateurs
- Probabilité, Histogramme, Densité, Fonction de répartition, Spectre de signal
- Ensembles, Partitions, Graphes dont arbres
- Génomes

Exemple $X = [0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 1]$,

$Y = [1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1]$

$(n_{10} + n_{01}) / (n_{00} + n_{11})$

Discordances/ Concordances



Semblable? quand?

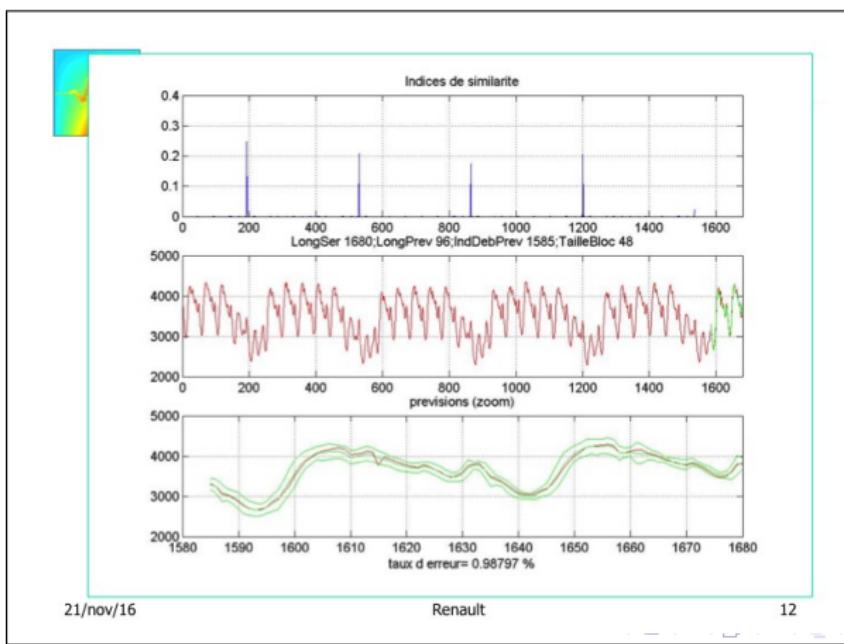
Deux éléments sont semblables si l'indice de dissimilarité est petit
 A partir d'un indice et d'un seuil

- De façon experte
- Si on dispose de la loi de variation: Test d'égalité à 0
- $U = r\sqrt{[n - 2]/(1 - r^2)}$, U distribué comme t à $n - 2$ ddl
- $V = (1/2)\log([1 + r]/[1 - r])$ distribué comme $N(0, 1/\sqrt{n - 3})$

1 Prévoir par similarité les analogues. Un exemple

- Indices de similarité un gros catalogue
 - distances
- Similarité de partitions
- Entropie
- Multidimensional Scaling
- Apprentissage de variété: distances géodésiques. ISOMAP.
- corrélations 2 a 2

Prévision par similarité 1



Prévision par similarité 2

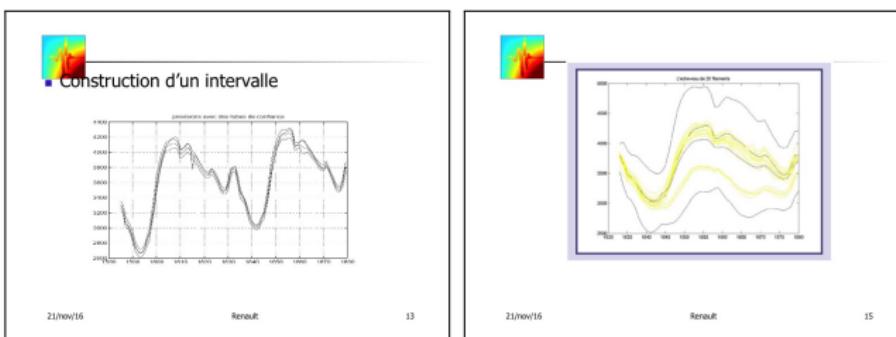


Figure: Que se passe t-il dans l'intervalle de confiance?

Distance



■ Distances

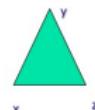
- $d(x,y)=0 \iff x=y$

- Symétrique

- Transitif

$$d(x,y) = d(y,x)$$

$$d(x,z) \leq d(x,y) + d(y,z)$$



22

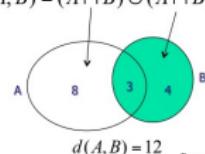


Similarité des graphes, une longue histoire

■ Distances entre ensembles (puis relations)

- Une vraie distance le nombre d'éléments de la différence symétrique

$$\Delta(A, B) = (A \cap \bar{B}) \cup (\bar{A} \cap B); d(A, B) = \#\Delta(A, B)$$



21/nov/16

$$d(A, B) = 12$$

Renault

23

Distance

Distance Euclidienne

- Il y a un produit scalaire $\langle x, y \rangle$

$\langle x, y \rangle$ forme bilinéaire positive

$\|x\|$ longueur $|x| = \sqrt{\langle x, x \rangle}$

$d(x, y)$ distance $\|x - y\|$

Perpendiculaires
Notion d'angle
(donc corrélation)

24

**Perpendiculaire, Pythagore
angle, cos
corrélation**

| Espace euclidien | Espace normé | Espace métrique | Espace topologique |
|------------------|--------------|-----------------|--------------------|
| ps norme | norme | | |
| distance | distance | distance | |
| topologie | topologie | topologie | topologie |

- Utilité: mesurer, proche ou lointain
- Structurer l'espace: boules
- Convergence: proximité en mouvement d'une position

25

Distance

 Ω ouvert de R^p

$H^1(\Omega) = \{f \in L^2(\Omega) / \frac{\partial f}{\partial x_j} \in L^2(\Omega)\}$

ps

$\langle f, g \rangle = \int_{(\Omega)} \left[(f(x)g(x) + \sum_{j=1}^p \frac{\partial f}{\partial x_j}(x) \frac{\partial g}{\partial x_j}(x) \right] d\mu(x)$

$H^1(\Omega)$ muni du produit scalaire est un espace de Hilbert

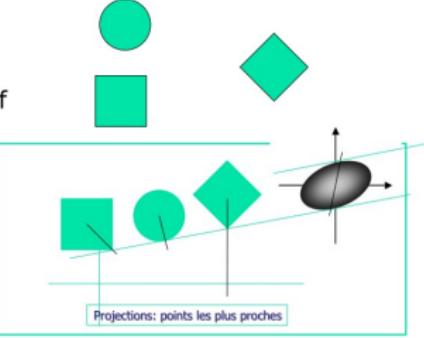
Si $p > 1$, les fonctions ne sont pas nécessairement continues mais on a un théorème disant qu'elles sont approchables, d'autant près que souhaité, par des fonctions très régulières

Ondelettes bien adaptées pour mesurer les écarts

27

 Métriques Minkovski (p)

- L2
- L1
- Linf



28

21/nov/16



Variation totale

La notion (norme sur un bon espace) est utile dans l'analyse d'image et des fonctions

$$\begin{aligned}
 VT(P, Q) &= \sup_{A \in \text{evenement}} |P(A) - Q(A)| \\
 &= (1/2) \int |p(x) - q(x)| dx \\
 &= (1/2) \sum_x |P(x) - Q(x)| dx
 \end{aligned}$$

Apparait en applications dans des domaines aussi variés que le traitement d'image, le contrôle optimal, le calcul des variations, et plus largement l'analyse numérique

Entropie

Entropie

40

Entropie (J. C. Hertz Information Theory, Inference and Learning Algorithms, Cambridge University Press, 2003)

- Entropie et théorème associé
 $H(M) = -\sum f_i \log_2 P_i \quad (\int dP \log_2 dP) \quad \text{trig} \theta = 0$
- Choisir une lettre au hasard en anglais apporte une information de l'ordre de 4.1 bits. La lettre z apporte 10.4 bits; j: 10.7; et 3.5 bits; et 4.1

$$H(X) = 0 \quad (\bar{H}(X) = 0 \Leftrightarrow P_i = 1, \forall i)$$

$$(i.e., H_{\text{max}}(H(X)) = \log_2(1/(0.01)(0.01)(0.01)(0.01))) \Rightarrow \bar{E} = c \theta = \frac{1}{|A|} \cdot |A|$$

Redundancy(X) = $1 - \frac{H(X)}{\log_2(|A|)}$

41

ENTROPIE MUTUELLE (et pseudo DISTANCE)

Entropie mutuelle, relative, ou divergence de Kullback-Leibler

$$(C, A) \text{ dans } \text{Kull}(P, Q) = \sum Q(x) \frac{P(x)}{Q(x)} \quad \left(\int dP(x) \log_2 \frac{dP(x)}{dQ(x)} \right)$$

$$\text{DKull}(P, Q) = 0 \quad (i.e. m = P = Q)$$

Non symétrique ($\text{DKull}(P, Q) \neq \text{DKull}(Q, P)$)
appelle parties distinctes!

Utilis. en théorie de l'information
Reconnaissance de forme
Réseaux de neurones
Séparation de sources. Analyse en composantes indépendantes

42

Entropie conjointe, conditionnelle. Information mutuelle

$$H(X, Y) = \sum_{i,j} P(i,j) \log_2 \frac{1}{P(i,j)}$$

$$X, Y \text{ indépendants } (P(i,j) = P(i)P(j)) \quad H(X, Y) = H(X) + H(Y)$$

$$\text{Entropie Conditionnelle de } X \text{ sachant } Y: j \quad H(X|Y=j) = \sum_i P(i|j) \log_2 \frac{1}{P(i|j)}$$

$$\text{Entropie Conditionnelle de } X \text{ sachant } Y: H(X|Y) = \sum_j P(j) H(X|Y=j)$$

Information mutuelle:
Indique la réduction moyenne de l'entropie due à la connaissance d'une valeur de (X, Y) : $I(X;Y) = H(X) - H(X|Y)$. $I(X;Y) \geq 0$, symétrique.
Distance ENTROPIQUE: différence entre l'entropie conjointe et l'entropie conditionnelle sur une variable diminue: $D(X, Y) = H(X, Y) - I(X, Y)$

43

Similarités. Ressemblances

- Après recodage en distribution. Ressemblance de 2 probabilités
 $d^2 = \sum \frac{(p_i^1 - p_i^2)^2}{n^2}$
- Entropie mutuelle
 $H(P, Q) = \sum Q(j) \log_2 \frac{P(j)}{Q(j)} \quad \left(\int dP \log_2 \frac{dP}{dQ} \right)$

Séparation de sources. Analyse en composantes indépendantes

- Entre des ordres: $x(i) > x(j)$ dans les deux. Wilcoxon.

44

| Densité | Clustering ooo | Clustering. Images oooo | ACP et la famille oooooooo oooooooooooo oooooooooooo ooooo oooooo | ICA ooo | MDS oooooooo | Ressembler | TP3 |
|---------|-------------------|----------------------------|--|------------|-----------------|------------|-----|
|---------|-------------------|----------------------------|--|------------|-----------------|------------|-----|

Plan

1 Densité

2 Clustering

3 Clustering. Images

4 ACP et la famille

5 ICA

6 MDS

Les données sont des échantillons de réaction de cellules dont on veut étudier les évolutions dans le temps. Les cellules sont différentes aux différents instants et leur nombre n_t change avec le temps.

- Réfléchir à la façon de définir la qualité de résultats dans les situations non supervisées.
- Réfléchir à la façon d'expliquer les classes et donner les raisons pour lesquelles des individus se retrouvent dans une même classe.
- Comparer des classifications.
- Evaluer l'importance es individus par suppression.
- Construire des classes de cellules.

Bibliography

- [1] Breiman, L., Friedman, J., Stone, C.J. and Olshen, R.A. (1984) *Classification and regression trees*. Editions Chapman & Hall/CRC, Monterey, CA.
- [2] Chavent, M., Lechevallier, Y. and Briant, O. (2007) *DIVCLUS-T: a monotonic divisive hierarchical clustering method*. Computational statistics and data analysis. 52(2), 687-701.
- [3] Corter, J. E. and Gluck, M. A. (1992). *Explaining basic categories: Feature predictability and information*. Psychological Bulletin, 111(2), 291-303.
- [4] Ester, M., Kriegel, H. P., Sander, J., and Xu, X. (1996) *A density-based algorithm for discovering clusters in large spatial databases with noise*. Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, 226–231.
- [5] Fisher, D. H. (1987) *Knowledge Acquisition Via Incremental Conceptual Clustering*. Machine Learning 2, 2:139-72.
doi:10.1023/A:1022852608280.
- [6] Fraiman, R., Ghattas, B. and Svarc, M. (2013) *Interpretable clustering using unsupervised binary trees*. Advances in data analysis and classification, 7, 125–145.
- [7] Ghattas, B., Svarc, M. and Fraiman, R. (2013) *R-package for interpretable clustering using binary trees*. <http://lumimath.univ-mrs.fr/ghattas/CUBT.html>.
- [8] Gluck, M. A. and Corter, J. E. (1985). *Information, uncertainty and the utility of categories*. In Proceedings of the Seventh Annual Conference of the Cognitive Science Society. Hillsdale NJ: Lawrence Erlbaum.

Bibliography...

- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., and Witten, I. H. (2009). *The WEKA Data Mining Software: An Update*. SIGKD D Explorations, Volume 11, Issue 1.
- Hennig C. (2014) *fpc: Flexible procedures for clustering*. R package version 2.1-9, <http://CRAN.R-project.org/package=fpc>.
- Huang, Z. (1998) *Extensions to the k-modes algorithm for clustering large data sets with categorical values*. Data Mining and Knowledge Discovery, 2(3), 283–304.
- Hubert, L. and Arabie, P. (1985) *Comparing partitions*. Journal of Classification, 2(1), 193–218.
- Kodinriki, T. M. and Makwana, P. R. (2013) *Partitioning Clustering algorithms for handling numerical and categorical data: a review*, arXiv, [arXiv:1311.7219v1](https://arxiv.org/abs/1311.7219v1).
- Leisch, F. (2006) *A tool box for K-centroids cluster analysis*. Computational Statistics and Data Analysis, 51(2), 526–544.
- Liu, B., Xu, Y., and Yu, P. S. (2000) *Clustering Through Decision Tree Construction*. In Proceedings of the Ninth International Conference on Information and Knowledge Management, 20–29. CIKM ’00. New York, NY, USA: ACM. doi:10.1145/354756.354775.
- MacQueen, J. (1967) *Some methods for classification and analysis of multivariate observations*. Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Editions L. M. Le Cam & J. Neyman, 1, 281–297.
- Murtagh, F. (1985) *Multidimensional Clustering Algorithms*. COMPSTAT Lectures 4, Physica-Verlag, Vienna.
- Rokach, L. (2010) *A survey of Clustering Algorithms* In Maimon, O. and Rokach, L. (eds.). *Data Mining and Knowledge Discovery Handbook*, 269–298, Springer US.
- Weihs, C., Ligges, U., Luebke, K. and Raabe, N. (2005) *klaR Analyzing German Business Cycles*. In Baier, D., Decker, R. and Schmidt-Thieme, L. (eds.). *Data Analysis and Decision Support*, 335–343, Springer-Verlag, Berlin.