

# Nonhierarchical clustering methods

**Mohamed Nadif**

LIPADE, Université Paris Descartes, France

# Outline

## 1 Introduction

- Cluster Analysis

## 2 *k*-means algorithm

- *k*-means: a family of methods
- Principal points to be retained

## 3 Spherical *k*-means

## 4 *k*-means for contingency table

## 5 *k*-means for categorical data

## 6 Miscellaneous clustering methods

- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy *k*-means

## 7 Conclusion

## Clustering

- Aim: It seeks to obtain a reduced representation of the initial data
- Organization of data into homogeneous subsets "clusters" or "classes"

## Structure of clustering

- It can take different forms: partitions, sequence of encased partitions or hierarchical, overlapping clusters, clusters with high density, fuzzy clusters.

## Characteristics of these methods

- This course is devoted to the partitioning methods or nonhierarchical clustering

# Outline

## 1 Introduction

- Cluster Analysis

## 2 *k*-means algorithm

- *k*-means: a family of methods
- Principal points to be retained

## 3 Spherical *k*-means

## 4 *k*-means for contingency table

## 5 *k*-means for categorical data

## 6 Miscellaneous clustering methods

- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy *k*-means

## 7 Conclusion

## Description

- We keep the previous notation and we begin by describing the well-know  $k$ -means when the set to classify  $\Omega$  is measured by  $p$  continuous variables
- To look for the optimal partition  $\mathbf{z}$  it suffices to minimize the within-cluster variance  $W(\mathbf{z})$

$$W(\mathbf{z}) = \sum_{k=1}^g \sum_{i \in \mathbf{z}_k} \|\mathbf{x}_i - \bar{\mathbf{x}}_{\mathbf{z}_k}\|^2 = \sum_{k=1}^g \sum_{i=1}^n z_{ik} \|\mathbf{x}_i - \bar{\mathbf{x}}_{\mathbf{z}_k}\|^2.$$

which is equivalent to maximize the between-cluster variance

$$B(\mathbf{z}) = \sum_{k=1}^g \pi_k \|\bar{\mathbf{x}}_{\mathbf{z}_k} - \bar{\mathbf{x}}\|^2,$$

where  $\pi_k$  is the weight of the cluster  $\mathbf{z}_k$  and  $\bar{\mathbf{x}}$  is the vector center of all data. This equivalence is due to the decomposition of the total variance  $I$  of data

$$I = \sum_{i=1}^n \pi_i \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 = W(\mathbf{z}) + B(\mathbf{z})$$

## Description of $k$ -means

- The one-parameter optimization  $W(\mathbf{z})$  is equivalent to the optimization of the two-parameter optimization  $W(\mathbf{z}, \boldsymbol{\mu})$  (Discrete sum-of-squares (SSQ))

$$W(\mathbf{z}, \boldsymbol{\mu}) = \sum_{k=1}^g \sum_{i \in z_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 = \sum_{i=1}^n \sum_{k=1}^K z_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \quad (1)$$

where  $z_{ik} \in \{0, 1\}$  and  $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g)$  with  $\boldsymbol{\mu}_k$  from  $\mathbb{R}^p$  represents the center or prototype of the cluster  $z_k$ .

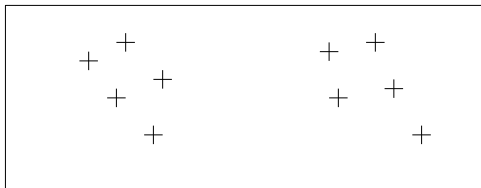
- This optimization can be carried out by the  $k$ -means algorithm and the principal steps of the  $k$ -means are the following:
  - 1 Randomly select  $g$  objects of  $\Omega$  which form the  $g$  first cluster means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g$ .
  - 2 While not convergence
    - 1 assign each object of  $\Omega$  to the cluster with the nearest cluster mean. If this one is not unique the object is assigned to the cluster with the smallest subscript.
    - 2 The cluster means computed become the new cluster means.

In the iteration process  $k$ -means yields a sequence  $\boldsymbol{\mu}^{(0)}, \mathbf{z}^{(1)}, \boldsymbol{\mu}^{(1)}, \mathbf{z}^{(2)}, \dots$  of partitions and centers with decreasing the values of the criterion until the convergence at the minimum value

## Description of $k$ -means

- We illustrate the different steps of  $k$ -means by applying it with  $g = 2$  on a simple set  $\Omega$  of 10 objects located in plan as depicted in a rectangle

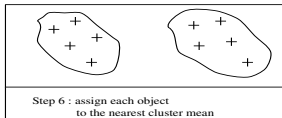
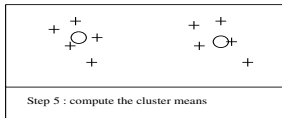
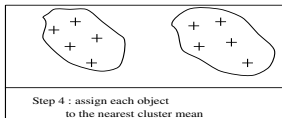
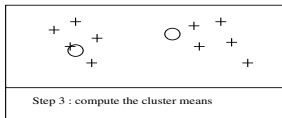
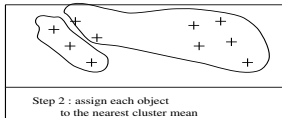
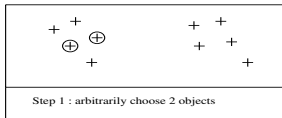
### Example of 10 objects to classify



## Description of $k$ -means (Forgy, 1965)

- The  $k$ -means algorithm can then be summarized in the following way

### Process of $k$ -means.



The process terminates and this algorithm will not change any more the results: The algorithm converges. Note that the obtained partition corresponds to the observable structure in two clusters



## Example

Id	x	y
S1	5	5
S2	6	6
S3	15	14
S4	16	15
S5	25	20
S6	26	19

## Without any calculation

- Initialize the  $k$ -means algorithm with S1, S4 et S6 and looking for 3 clusters.  
Remark ?
- Initialize the  $k$ -means algorithm with S4, S5 et S6 and looking for 3 clusters.  
Remark ?

## Within-cluster criterion

- It corresponds to the famous sum-of-squares criterion (SSQ)
- Different approaches in clustering are based on this criterion but under different forms due to different hypotheses

## Classical hypothesis

- $\mathbf{z}$  is a known partition and  $x_1, \dots, x_n$  as a realization of a random vector  $\mathbf{X}$  with  $f$  its density on  $\mathbb{R}^p$ .
  - The problem is to look for the partition  $\mathbf{z}$  in  $\mathbb{R}^p$  minimizing :

$$W(\mathbf{z}) = \sum_k \int_{z_k} \|\mathbf{x} - \mathbb{E}_{z_k}(X)\|^2 dP(\mathbf{x})$$

As  $P(x < X < x + dx) = f(x)dx$ ,  $W(\mathbf{z})$  is equivalent to

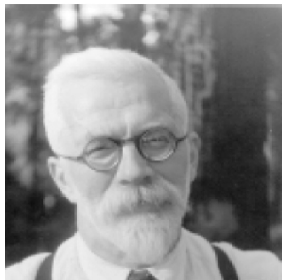
$$\begin{aligned} W(\mathbf{z}, \mu) &= \sum_k \int_{z_k} \|\mathbf{x} - \mathbb{E}_{z_k}(X)\|^2 f(\mathbf{x}) d\mathbf{x} \\ &= \sum_k \int_{z_k} f(\mathbf{x}) \|\mathbf{x} - \mu_k\|^2 d\mathbf{x} \end{aligned} \quad (2)$$

- In  $\mathbb{R}$ , this formulation has been provided in the framework of optimum proportional stratified sampling (Dalenius, 1950). Even if the  $k$ -means algorithm was not used but another called the *shooting* algorithm, this later needs two principal steps of  $k$ -means

## Within-cluster criterion

- Different extensions to this algorithm were proposed and successfully applied in image compression. The figure( photo-Fisher) illustrates an application of the LLoyd algorithm in the context of scalar quantization.

## Example of scalar quantization.



## A word about the within-cluster criterion

- the optimization of SSQ was considered in multidimensional case  $\mathbb{R}^p$ , and the first to propose the  $k$ -means explicitly was Steinhauss (1956)
- Actually, the  $k$ -means version commonly used is due to Forgy (1965)
- Another approach which consists to consider the data as a sample is appeared with MacQueen (1967) and a stochastic version of  $k$ -means is performed and inspired the Kohonen maps
- We can cite the dynamic clustering method (Diday, 1971), iterated minimum-distance partition method (Bock, 1974)
- After several works concerned on different variants of  $k$ -means and sometimes under different names algorithms were proposed, see for instance (Bock, 2007)

- Before performing a cluster analysis on coordinate data, it is necessary to consider scaling or transforming the variables
- If our aim is to find the couple  $(z, \mu)$  minimizing the criterion  $W$ , the  $k$ -means algorithm does not provide necessarily the best result, but just a sequence of couples whose the value of criterion is going to decrease and we obtain a local optimum.
- The  $k$ -means algorithm can use an  $L_m$  clustering criterion instead of the least-squares  $L_2$  criterion. Note that values of  $m$  less than 2 reduce the effect of outliers on the cluster centers compared with least-squares criterion
- In general, the criterion is not independent of the number of classes. For example, the partition into  $n$  classes, where each object forms a singleton cluster, has a null within-cluster criterion and therefore the optimal partition is without interest. It is then necessary to fix a priori the number of classes.

- If this number is not known, several solutions allowing to solve this very difficult problem are used. For example, the best partition is sought for several numbers of classes and we study the decrease of the criterion according the number of classes to select the number of classes by using the scree plot and choosing an elbow. Indeed, the quality of a partition can be evaluated by the *Rsquare* (RSQ)

$$RSQ = 1 - \frac{W}{I} = \frac{B}{I}$$

The *k*-means has a very low computational complexity which translates directly into a high speed, it suffices then to run *k*-means with different number of clusters and use the elbow method. Different criteria are available in NbClust.

- Knowing that according to starting points chosen, the results will be different, it remains with to exploit these different results. Several solutions were proposed: we run the *k*-means several times by initiating with different random initializations. Several strategies are then possible.
  - We select a *good* initialization with supplementary informations or with an automatic procedure (points strongly distant, regions with high density, etc.).
  - We should however make a compromise between the necessary time to the research of the initial configuration and that necessary time for the algorithm itself.

## Hybrid method (Wong, 1982)

- Link between  $k$ -means and the Ward method: The two methods are similar in that they both attempt to the within-cluster variance.
  - 1 Apply  $k$ -means to cluster  $\Omega$  into fifty clusters, for example. In practice, this number depending on the size of data can be taken equal to  $n^{\frac{1}{3}}$
  - 2 Run the Ward method on these obtained cluster means
  - 3 From the dendrogram we propose a number of clusters by using the SPRSQ criterion
  - 4 Eventually, apply  $k$ -means on the obtained clusters to improve SPRSQ
- Fisher's method (1958): Note that there exist some situations for which we have effective algorithms allowing to find global optimum. It is the case where there is an order constraint on the partitions. This constraint can be implicit (for instance, when the data are in  $\mathbb{R}$ ) or explicit (for instance, constraint imposed by the user). We can then use a dynamic algorithm of programming such as the Fisher's algorithm which provides the global optimum
- Interpretation of clusters (exploratory tools, ANOVA, etc.)

# Outline

## 1 Introduction

- Cluster Analysis

## 2 $k$ -means algorithm

- $k$ -means: a family of methods
- Principal points to be retained

## 3 Spherical k-means

## 4 $k$ -means for contingency table

## 5 $k$ -means for categorical data

## 6 Miscellaneous clustering methods

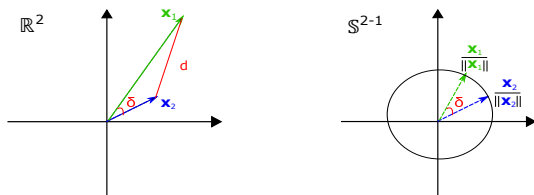
- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy  $k$ -means

## 7 Conclusion



## Challenges

- The high dimensionality and sparsity characterising the data sets arising in some areas, such as Recommender systems and text mining
- Such data sets consist of more than 1000 features and 95% of zero entries
- The data sets from the aforementioned domains are also directional in nature



- Popular clustering based on the euclidean distance, for instance, are inadequate for directional data

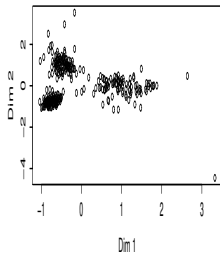
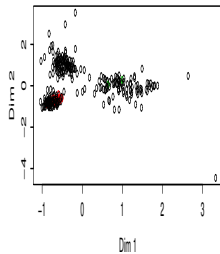
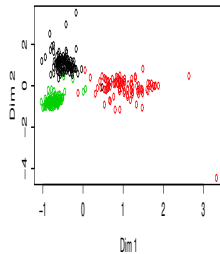
## Remarks

- The spherical k-means is tailored for directional data distributed on the surface of a unit-hypersphere
- Text document clustering, microarray-data, and item recommendation are the popular domains where this algorithm is effective.

## Spherical k-means: classic300 ( $300 \times 5577$ )

```
library(R.matlab)
setwd("/Users/nadif/Desktop/China/docdata")
classic300 <- readMat("classic300.mat")
dim(classic300$dtm)
x=as.matrix(classic300$dtm)
#"Correspondence Analysis"
library(FactoMineR)
res.ca=CA(x)
plot(res.ca,invisible="col")
#"application of kmeans and visualisation of clusters"
z.kmeans <- kmeans(x, 3, nstart = 100)
table(z.kmeans$cluster,classic300$classid)
plot(res$row$coord,col=z.kmeans$cluster)
#"application of skmeans and visualisation of clusters"
library(skmeans)
zs.skmeans <- skmeans(x, 3)
table(zs.skmeans$cluster,classic300$classid)
plot(res$row$coord,col=zs.skmeans$cluster)
```

Original data

 $k$ -meansSpherical  $k$ -means

# Outline

## 1 Introduction

- Cluster Analysis

## 2 $k$ -means algorithm

- $k$ -means: a family of methods
- Principal points to be retained

## 3 Spherical $k$ -means

## 4 $k$ -means for contingency table

## 5 $k$ -means for categorical data

## 6 Miscellaneous clustering methods

- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy  $k$ -means

## 7 Conclusion

## Clustering of contingency table

- The classical Euclidean distance is not appropriated
- Clusters of objects is not informative than clustering of the row profiles
- The  $\chi^2$  is more adapted

## Transformation of data

- The data matrix  $\mathbf{X}$  is a  $n \times p$  matrix defined by  $\mathbf{X} = \{(x_{ij}); i \in I, j \in J\}$  where  $I$  is a categorical variable with  $n$  categories and  $J$  a categorical variable with  $p$  categories
- We denote the row and columns total of  $\mathbf{x}$  by  $x_{i.} = \sum_{j=1}^p x_{ij}$  and  $x_{.j} = \sum_{i=1}^n x_{ij}$  and the overall total simply by  $N = \sum_{i,j} x_{ij}$ .
- We denote  $\{(f_{ij} = x_{ij}/N); i \in I, j \in J\}$
- the marginal frequencies  $f_{i.} = \sum_j f_{ij} = x_{i.}/N$  and  $f_{.j} = \sum_i f_{ij} = x_{.j}/N$
- The row profiles  $f_i^J = (f_{i1}/f_{i.}, \dots, f_{ip}/f_{i.})^T = (x_{i1}/x_{i.}, \dots, x_{ip}/x_{i.})^T$
- The average row profile  $f_J = (f_{.1}, \dots, f_{.p})^T = \frac{1}{N}(x_{.1}, \dots, x_{.p})^T$

## Contingency table and $\chi^2$

There are several measures of association and the most employed is the chi-square  $\chi^2$ . This criterion, used for example in the correspondence analysis

$$\chi^2(I, J) = \sum_{i=1}^n \sum_{j=1}^p \frac{(x_{ij} - \frac{x_{i.} x_{.j}}{N})^2}{\frac{x_{i.} x_{.j}}{N}} = N \sum_{i=1}^n \sum_{j=1}^p \frac{(f_{ij} - f_{i.} f_{.j})^2}{f_{i.} f_{.j}}.$$

- $\chi^2$  usually provides statistical evidence of a significant association, or dependence, between rows and columns of the table. It represents the deviation between the theoretical frequencies  $f_{i.} f_{.j}$ , that we would have if  $I$  and  $J$  were independent, and the observed frequencies  $f_{ij}$ . If  $I$  and  $J$  are independent, the  $\chi^2$  will be zero and if there is a strong relationship between  $I$  and  $J$ , the  $\chi^2$  will be high
- A significant chi-square indicates a departure from row or column homogeneity and can be used as a measure of heterogeneity. Then, the chi-square can be used to evaluate the quality of a partitions of  $I$  or  $\mathbf{w}$  of  $J$
- Associated  $\chi^2(\mathbf{z}, J)$  of the contingency table with  $K$  rows in making the sum of rows of each cluster. We have  $\chi^2(I, J) \geq \chi^2(\mathbf{z}, J)$  and the objective is to find the partitions  $\mathbf{z}$  which minimizing this loss, i.e. which maximizes

$$\chi^2(\mathbf{z}, J) = N \sum_{k=1}^K \sum_{j=1}^p \frac{(f_{kj} - f_{k.} f_{.j})^2}{f_{k.} f_{.j}}$$

## Time-budget data matrix

	prof	tran	home	child	shop	wash	meal	sleep	tv	leis
maus	610	140	60	10	120	95	115	760	175	315
waus	475	90	250	30	140	120	100	775	115	305
wnaus	10	0	495	110	170	110	130	785	160	430
mnsus	615	141	65	10	115	90	115	765	180	305
wnsus	179	29	421	87	161	112	119	776	143	373
msus	585	115	50	0	150	105	100	760	150	385
...	...	...	...	...	...	...	...	...	...	...
...	...	...	...	...	...	...	...	...	...	...
mnsea	652	133	134	22	68	94	102	762	122	310
wnsea	434	77	431	60	117	88	105	770	73	229
msea	627	148	68	0	88	92	86	770	58	463
wsea	433	86	296	21	128	102	94	758	58	379

$I$  : types of population and  $J$  : variety of activities,  $x_{ij}$ : amount of time spent on a variety of activities  $j$  by  $i$  during a given time period  $j$

## Notation

- The choice of  $\chi^2$  metric is justified for several reasons, in particular because of the similar role played by each of the two dimensions in the analyzed table, and also because of the property of distributional equivalence, which implies stable results when agglomerating elements with similar profiles
- Each row  $i$  corresponds to a point vector  $\mathbb{R}^p$  defined by the profile  $f_{iJ}$  weighted by the marginal frequency  $f_i$ .
- The maximization of  $\chi^2(\mathbf{z}, J)$  can be viewed as the minimization of a criterion depending on the partition and the centers of clusters
- $\mathbf{z}$  is a partition of the rows, we can define the frequencies  $f_{kj} = \sum_{i,k} z_{ik} f_{ij}$  and the average row profile of the  $k$ th cluster is defined by the vector  $f_{kJ} = (\frac{f_{k1}}{f_{k.}}, \dots, \frac{f_{kp}}{f_{k.}})^T$  where  $f_{k.} = \sum_{j=1}^p f_{kj}$



## SSQ criterion

- With this representation, the total of squared distances  $T$ , the between-cluster sums of squares  $B(\mathbf{z})$  and the within-cluster sums of squares take the forms

$$T = \sum_{i=1}^n f_i. d^2(f_{iJ}, f_J) , \quad B(\mathbf{z}) = \sum_{i=1}^n f_k. d^2(f_{kJ}, f_J) = \frac{1}{N} \chi^2(\mathbf{z}, J),$$

and

$$W(\mathbf{z}) = \sum_{k=1}^g \sum_{i=1}^n z_{ik} f_i. d^2(f_{iJ}, f_J).$$

- The traditional relation  $T = W(\mathbf{z}) + B(\mathbf{z})$  leads to the following relation:

$$\chi^2(I, J) = NW(\mathbf{z}) + \chi^2(\mathbf{z}, J).$$

- The term  $NW(\mathbf{z})$  therefore represents the information lost when grouping the elements according to the partition  $\mathbf{z}$ , and  $\chi^2(\mathbf{z}, J)$  corresponds to the information which is preserved. Looking for the partition maximizing the criterion  $\chi^2(\mathbf{z}, J)$  is equivalent to looking for the partition minimizing  $W(\mathbf{z})$  or  $W(\mathbf{z}, \mathbf{a})$
- To minimize this criterion it is possible to apply  $k$ -means to the set of profiles with the  $\chi^2$  metric. The iterative algorithm maximizes locally  $\chi^2(\mathbf{z}, J)$

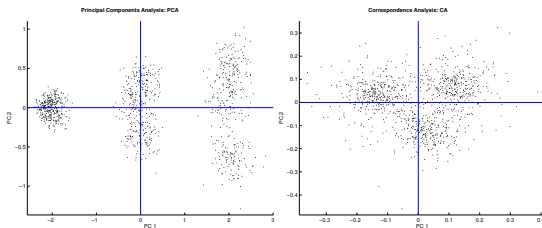
## Importance of the metric

- Let be  $\mathbf{X}$  a data matrix that consists of a set of objects described by 3 continuous variables  $x$ ,  $y$  and  $z$ . Naturally, we can use the Standardized Euclidean distance on standardized data but sometimes the clustering on the profiles (row percents) are more adapted in certain contexts and as the values of this data matrix are all positive we can use the metric  $\chi^2$

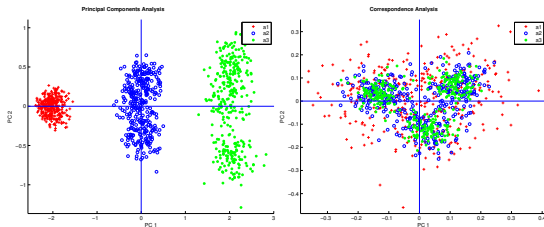
## Extract of data

x	y	z	x	y	z
37	31	40	117	132	142
35	26	29	166	118	117
42	44	25	115	126	153
43	20	28	152	105	115
32	26	43	114	119	162
44	32	27	109	109	91
31	38	29	136	150	95
28	47	49	100	132	152
..	..	..	...	...	...

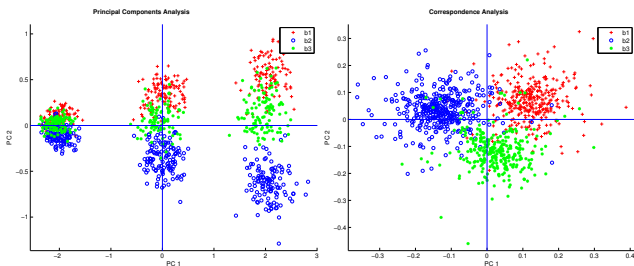
## Projection of objects on the factorial planes spawned by the first and second axes by PCA and CA



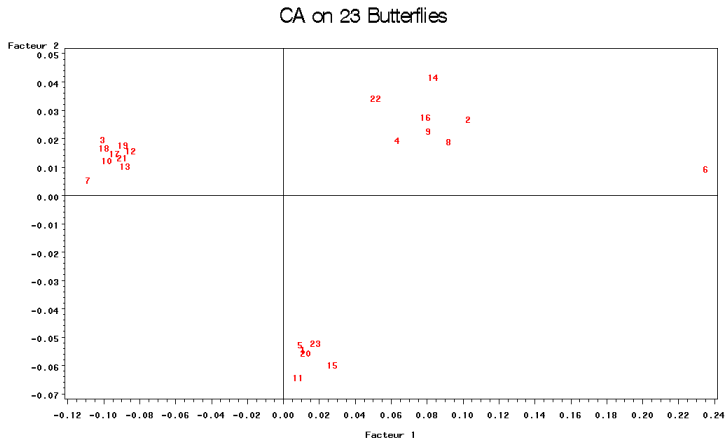
Projection of clusters (obtained with the standard euclidian distance) on the factorial planes spawned by the first and second axes by PCA and CA.



Projection of clusters (obtained with the  $\chi^2$  distance) on the factorial planes spanned by the first and second axes by PCA and CA.



## Example of Butterflies



### k-means by using the $\chi^2$ distance

- This table can be considered as the contingency table
- Why ?

# Outline

## 1 Introduction

- Cluster Analysis

## 2 *k*-means algorithm

- *k*-means: a family of methods
- Principal points to be retained

## 3 Spherical *k*-means

## 4 *k*-means for contingency table

## 5 *k*-means for categorical data

## 6 Miscellaneous clustering methods

- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy *k*-means

## 7 Conclusion

## Clustering of Categorical data

- Generally we apply the clustering to a particular indicator matrix
- Let a variable with 3 categories 1, 2, 3  $\Rightarrow (1, 0, 0), (0, 1, 0), (0, 0, 1)$ , then the matrix has the number of rows equal to the total number of objects and the number of columns equal to the sum of all categories corresponding to all variables
- As before the  $\chi^2$  is the more appropriate metric and we can apply the kmeans with the  $\chi^2$  metric

	a	b		a1	a2	a3	b1	b2	b3
1	1	2	1	1	0	0	0	1	0
2	3	2	2	0	0	1	0	1	0
3	2	3	3	0	1	0	0	0	1
4	1	1	4	1	0	0	1	0	0
5	1	2	5	1	0	0	0	1	0
6	3	2	6	0	0	1	0	1	0
7	3	3	7	0	0	1	0	0	1
8	1	1	8	1	0	0	1	0	0
9	2	2	9	0	1	0	0	1	0
10	2	3	10	0	1	0	0	0	1

The  $\chi^2$  distance takes the following form:

$$d_{\chi^2}(i, i') = \sum_{j=1}^p \frac{1}{f_{.j}} \left( \frac{f_{ij}}{f_{i.}} - \frac{f_{i'j}}{f_{i'.}} \right)^2 = \sum_{j=1}^p \frac{N}{x_{.j}} \left( \frac{x_{ij}}{x_{i.}} - \frac{x_{i'j}}{x_{i'.}} \right)^2 \propto \sum_{j=1}^p \frac{1}{x_{.j}} \left( \frac{x_{ij}}{x_{i.}} - \frac{x_{i'j}}{x_{i'.}} \right)^2$$

then  $d_{\chi^2}(i, i') \propto \sum_{j=1}^p \left( \frac{x_{ij}}{\sqrt{x_{.j}x_{i.}}} - \frac{x_{i'j}}{\sqrt{x_{.j}x_{i'.}}} \right)^2$

# Outline

## 1 Introduction

- Cluster Analysis

## 2 $k$ -means algorithm

- $k$ -means: a family of methods
- Principal points to be retained

## 3 Spherical $k$ -means

## 4 $k$ -means for contingency table

## 5 $k$ -means for categorical data

## 6 Miscellaneous clustering methods

- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy  $k$ -means

## 7 Conclusion



## Sequential methods

- The  $k$ -means algorithm has different extensions in order to apply it on different types of data such as the sequential data
- The online variants of  $k$ -means are particularly adequate when all the data to be classified are not available in the beginning
- The parameters defining the classes can then be adjusted when a new data comes as a continuous stream without too many calculations
- Unlike  $k$ -means the objects concerning by the step assignment are randomly selected and the update of cluster means is realized after each assignment of one object
- More precisely, at the  $(t)$ th iteration, the object  $x_i$  is randomly selected, then we determine the nearest prototype  $\mu_k^{(t)}$  which becomes after assignment of  $x_i$  equal to

$$\mu_k^{(t+1)} = \frac{x_i + n_k^{(t)} \cdot \mu_k^{(t)}}{n_k^{(t+1)}},$$

where  $n_k^{(t)}$  represents the cardinality of the cluster  $z_k^{(t)}$  and  $n_k^{(t+1)} = n_k^{(t)} + 1$ .

$$\mu_k^{(t+1)} = \mu_k^{(t)} + \frac{1}{n_k^{(t+1)}}(x_i - \mu_k^{(t)}),$$

## Self-Organizing-Mapping

- *Self-Organizing-Mapping* or SOM a type of clustering, inspired by neuroscience, that has been introduced in Kohonen (1982).
- In the SOM literature, we refer to the clusters by the nodes or neurons and each of them has a weight in  $\mathbb{R}^p$ , these weights refer to the cluster means
- The principal advantage of SOM that is preserves the topology clustering. Generally, the neurons are arranged as one or two-dimensional rectangular grid preserving relations between the objects called also units
- SOM offers a good tool to visualize clusters and evaluate their proximity in a reduced space
- Unlike  $k$ -means and AHC, the previous expression of the cluster means or the weight of a neuron  $k$  becomes in the SOM context

$$\mu_k^{(t+1)} = \mu_k^{(t)} + \varepsilon(t) \times h(k, \ell)(x_i - \mu_k^{(t)}), \quad (3)$$

where  $h(k, \ell)$  is the neighborhood function between the neuron  $k$  whose weight  $\mu_k^{(t)}$  is the most similar to  $x_i$  (the best matching unit or the winner) and the other neurons  $\ell$  with weights  $\mu_\ell$  close enough to  $\mu_k^{(t)}$

## The neighborhood function

- This function can take different forms, it evaluates the proximity between the winner  $k$  and the neuron  $\ell$  located in a reduced space generally in  $\mathbb{R}^2$  with the position  $r_k$  and  $r_\ell$ 
  - In the early publications about SOMs,  $h(k, \ell)$  was defined by:  $h(k, \ell) = 1$  if  $d(k, \ell) \leq \lambda$  and 0 otherwise
  - Gaussian function is a common choice  $h(k, \ell) = \exp(-\frac{\alpha \|r_k - r_\ell\|^2}{2\sigma_h(t)})$  where  $\sigma_h(t)$  controls the width of the neighborhood of  $h$ .

## SOM Algorithm

- A neuron  $k$  is characterized by the weight vector  $\mu_k$
- Description of the basic SOM
  - Choose the size of the grid initialization of the neurons:  $\mu_k^{(0)}$
  - Choose an object  $x_i^{(c+1)}$
  - Research of the winner  $k^*$ ,  $k^* = \operatorname{argmin}_k \|x_i^{(c+1)} - \mu_k^{(c)}\|$
  - The update of the weight vectors concern  $k^*$  and all neurons near of  $k^*$

$$\mu_k^{(c+1)} = \mu_k^{(c)} + \varepsilon(t)h(k^*, \ell)(x_i^{(c+1)} - \mu_k^{(c)})$$

## Characteristics of SOM

- With this grid moving during the iterations of SOM, we obtain a partition and a visualization of the clusters as in factorial plan from PCA, except this representation is not linear because it is not an orthogonal projection
- The different steps of SOM are similar than the steps of  $k$ -means. In addition, two versions batch and online can be used. The first one performs the assignment and update steps for all data units at once and the second process as the MacQueen algorithm
- As  $k$ -means, SOM requires to fix the number of clusters (nodes of the grid), and a choice of initialization. PCA appears an attractive and interesting approach
- As the numbers of nodes is higher, the number of clusters can be assess by applying AHC algorithm with appropriated agglomerative criterion on these nodes

- The mean can be generalized

$$\mu_k^{(t+1)} = \mu_k^{(t)} + \varepsilon(t)(x_i - \mu_k^{(t)}),$$

where  $\varepsilon(t)$  is a decreasing learning coefficient

- The usual hypothesis on the adaptation parameter to get almost sure results is then (conditions of Robbins-Monro):

$$\sum_t \varepsilon(t) = +\infty \text{ and } \sum_t \varepsilon(t)^2 < +\infty$$

- This formulation of the cluster means can be extended and constitutes the version of other algorithms such as the well-known *Self-Organizing-Mapping*

## LSA: Latent Semantic Analysis (Deerwester et al., 1988)

- LSA (or LSI) is a technique in NLP of analyzing relationships between the set of documents and the terms they contain by producing a set of concepts (components in PCA) related to the documents and terms
- LSA is quite similar to PCA
- Let  $\mathbf{X}$  be document-term matrix, LSA is based on SVD:  
 $\mathbf{X}_{(n \times p)} = \mathbf{U}_{(n \times p)} \mathbf{\Sigma}_{(p \times p)} \mathbf{V}_{(p \times p)}^T$  where  $\mathbf{U}\mathbf{U}^T = \mathbf{I}$  and  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$
- LSA finds a low-rank approximation  $\mathbf{X} \approx \mathbf{U}_{(n \times g)} \mathbf{\Sigma}_{(g \times g)} \mathbf{V}_{(g \times p)}^T$
- $\mathbf{XV} = \mathbf{U}\mathbf{\Sigma}$  is the set of terms and  $\mathbf{X}^T\mathbf{U} = \mathbf{V}\mathbf{\Sigma}$  is the set of documents in the latent topic

## Advantages and limitations

- + LSA allows a comparison of documents or terms in the low-dimensional space
  - The resulting dimensions might be difficult to interpret
  - LSA cannot capture polysemy (multiple meaning of a term, a query "Tree" will not give expected documents for a botanist or for a statistician)

## NMF: Nonnegative Matrix Factorization (Lee and Seung, 1999, 2001)

- Problem:  $\operatorname{argmin}_{\mathbf{U}, \mathbf{V} \geq 0} \|\mathbf{X} - \mathbf{UV}^T\|^2$  where factor matrices,  $\mathbf{U} \in \mathbb{R}_+^{n \times g}$  and  $\mathbf{V} \in \mathbb{R}_+^{p \times g}$
- l-divergence:  $\operatorname{argmin}_{\mathbf{U}, \mathbf{V} \geq 0} = \sum_{i=1}^n \sum_{j=1}^p X_{ij} \log \frac{X_{ij}}{(UV)_{ij}} - X_{ij} + (UV)_{ij}$
- The clustering problem is not the main objective of NMF
- Each column of  $\mathbf{X}$  is treated as a data point in  $n$ -dimensional space
- Each  $U_{ik}$  of  $\mathbf{U}$  corresponds to the degree to which row  $i$  belongs to  $k$ th cluster
- Each column of  $\mathbf{U}$  is associated with a prototype vector for the  $k$ th cluster

$$\mathbf{X} = \mathbf{U} \mathbf{V}^T$$

### Advantages and limitations

- + Unlike LSA, NMF does not require the derived latent semantic space to be orthogonal
- + Unlike LSA, the interpretation is more easy due to the positivity of values
  - Uniqueness, initialization

## Expressions of $\mathbf{U}$ and $\mathbf{V}$

A typical constrained optimization problem can be solved using the Lagrange multiplier method:  $U_{ik} \leftarrow U_{ik} \frac{(\mathbf{XV})_{ik}}{(\mathbf{UV}^T\mathbf{V})_{ik}}$  and  $V_{kj} \leftarrow V_{kj} \frac{(\mathbf{X}^T\mathbf{U})_{kj}}{(\mathbf{VU}^T\mathbf{U})_{kj}}$

## Uniqueness

If  $\mathbf{U}$  and  $\mathbf{V}$  are solutions, then  $\mathbf{UD}^{-1}$  and  $\mathbf{VD}$  will also form a solution for any positive diagonal matrix  $\mathbf{D}$ . To eliminate this uncertainty, in practice one will further require that the Euclidean length of each column vector in  $\mathbf{U}$  is 1  $U_{ik} \leftarrow \frac{U_{ik}}{\sqrt{\sum_i U_{ik}^2}}$ ,  $V_{kj} \leftarrow V_{kj} \sqrt{\sum_i U_{ik}^2}$

## Document clustering: NMF towards clustering (library(NMF))

- 1 Perform the NMF on  $\mathbf{X}$  to obtain  $\mathbf{U}$  and  $\mathbf{V}$
- 2 Normalize  $\mathbf{U}$  and  $\mathbf{V}$
- 3 Use matrix  $\tilde{\mathbf{U}}$  to determine the cluster label of each document. Examine each document of matrix  $\tilde{\mathbf{U}}$  and assign it to cluster  $k^*$  if  $k^* = \arg \max_k \tilde{\mathbf{U}}_{ik}$

## Orthogonal NMF

$\operatorname{argmin}_{\mathbf{U}, \mathbf{V} \geq 0} \|\mathbf{X} - \mathbf{UV}^T\|^2$  where  $\mathbf{U} \in \mathbb{R}_+^{n \times g}$ ,  $\mathbf{V} \in \mathbb{R}_+^{g \times p}$ ,  $\mathbf{UU}^T = \mathbf{I}_n$ ,  $\mathbf{V}^T\mathbf{V} = \mathbf{I}_p$



## Dynamic clusters method

- The dynamic clusters method (Diday, 1971). The cluster centers are not necessarily cluster means elements of  $\mathbb{R}^p$ .
- Let  $\mathbb{L}$  be the set of centers, and  $D : \Omega \times L \rightarrow \mathbb{R}^+$ , a measure of dissimilarity between objects of  $\Omega$  and the centers of  $\mathbb{L}$ . The aim is to look for a partition of  $\Omega$  into  $g$  clusters minimizing the following criterion

$$C(\mathbf{z}, L) = \sum_{k=1}^g \sum_{i=1}^n z_{ik} D(\mathbf{x}_i, \lambda_k)$$

where  $\mathbf{z} = (z_{ik})$  with  $z_{ik} \in \{0, 1\}$   $L = (\lambda_1, \dots, \lambda_K)$  with  $\lambda_k \in \mathbb{L}$

- If  $\Omega \subset \mathbb{R}^p$ ,  $\mathbb{L} = \mathbb{R}^p$  and  $D(\mathbf{x}, \lambda) = d^2(\mathbf{x}, \lambda)$  then  $C(\mathbf{z}, L) = W(\mathbf{z}, \mu)$ .
- Like  $k$ -means, to tackle the minimization of  $C(\mathbf{z}, L)$  we can use an alternating optimization method based on
  - 1 Compute  $\mathbf{z}^{(t+1)}$  minimizing  $C(\cdot, L^{(t)})$
  - 2 Compute  $L^{(t+1)}$  minimizing  $C(\mathbf{z}^{(t+1)}, \cdot)$ .
- These steps yield the following sequence

$$L^{(0)} \rightarrow \mathbf{z}^{(1)} \rightarrow L^{(1)} \rightarrow \mathbf{z}^{(2)} \rightarrow L^{(2)} \rightarrow \dots \rightarrow \mathbf{z}^{(t)} \rightarrow L^{(t)} \rightarrow \dots$$

where  $L^{(0)}$  is an arbitrary initialization

The dynamical clusters method became so classical that often it is referred wrongly by  $k$ -means. It allows the user to choose the nature of the cluster centers. Next, we sketch different situations

### Examples of dynamic clusters method

- The  $k$ -medoids algorithm is a typical dynamical cluster methods minimizing the SSQ criterion, But in contrast to  $k$ -means the cluster centers are objects of  $\Omega$ . The principal advantage of  $k$ -medoids is that it overcomes the problem of outliers (Kaufman, 1987). A version PAM (Partition Around Medoids) by Kaufman (1990). Other variants of PAM with less complexity were proposed such as CLARANS by Ng and Han (1994)
- The choice of distance is crucial in clustering. In the process of dynamical clusters we can try to learn a metric. Then instead of to use a fixed distance we can consider that  $\Omega \subset \mathbb{R}^p$  and  $\mathbb{L} = \mathbb{R}^p \times \Delta$  where  $\Delta$  is set of distances defined on  $\mathbb{R}^p$  and  $D(x, (\lambda, d)) = d(x, \lambda)$ . Then the method performs clustering and distance metric learning simultaneously. This process allows to take into account the shapes of clusters (Diday, 1974, 1977)

## Clustering of categorical data

- The dissimilarity between two vectors of categories can be expressed as

$$D(\mathbf{x}_i, \boldsymbol{\lambda}_k) = \sum_{j=1}^p \delta(\mathbf{x}_{ij}, \boldsymbol{\lambda}_{kj})$$

where  $\delta(\mathbf{x}_{ij}, \boldsymbol{\lambda}_{kj}) = 1$  if  $\mathbf{x}_{ij} = \boldsymbol{\lambda}_{kj}$  and 0 otherwise

- $D$  reflects the number of different categories between the vector  $\mathbf{x}_i$  and the center  $\boldsymbol{\lambda}_k$   
Let  $\mathbf{x}_i = (1, 3, 2, 1, 2, 3)^T$  and  $\mathbf{x}_{i'} = (1, 1, 3, 1, 2, 1)^T$

$$D(\mathbf{x}_i, \mathbf{x}_{i'}) = 0 + 1 + 1 + 0 + 0 + 1 = 3$$

- The centers of obtained clusters are summarized by the vectors of categories (Marchetti, 1991) or ( $k$ -modes by Huang, 97)

## Binary data

- When the data are binary, the distance

$$D(\mathbf{x}, \boldsymbol{\lambda}_k) = \sum_{j=1}^p |\mathbf{x}_{ij} - \boldsymbol{\lambda}_{kj}|$$

is the Manhattan distance and the vector centers belong to  $\{0, 1\}^p$

## Nominal categorical data matrix and reorganized data matrix

	a	b	c	d	e		a	b	c	d	e
1	1	2	2	3	2	3	2	3	3	1	1
2	3	2	1	1	1	7	3	3	2	1	1
3	2	3	3	1	1	9	2	2	2	1	1
4	1	1	2	3	3	10	2	3	3	2	2
5	1	2	1	3	3	1	1	2	2	3	2
6	3	2	1	1	2	4	1	1	2	3	3
7	3	3	2	1	1	5	1	2	1	3	3
8	1	1	1	3	3	8	1	1	1	3	3
9	2	2	2	1	1	2	3	2	1	1	1
10	2	3	3	2	2	6	3	2	1	1	2

## Centers and Degree of homogeneity

	a	b	c	d	e		a	b	c	d	e
A	2	3	2	1	1	A	75	75	50	75	75
B	1	1	1	3	3	B	100	50	50	100	75
C	3	2	1	1	1	C	100	100	100	100	50

## Binary data matrix and reorganized data matrix

	a	b	c	d	e		a	b	c	d	e
1	1	0	1	0	1	1	1	0	1	0	1
2	0	1	0	1	0	4	1	0	1	0	0
3	1	0	0	0	0	8	1	0	1	0	1
4	1	0	1	0	0	2	0	1	0	1	0
5	0	1	0	1	1	5	0	1	0	1	1
6	0	1	0	0	1	6	0	1	0	0	1
7	0	1	0	0	0	10	0	1	0	1	0
8	1	0	1	0	1	3	1	0	0	0	0
9	1	0	0	1	0	7	0	1	0	0	0
10	0	1	0	1	0	9	1	0	0	1	0

## Centers and Degree of homogeneity

	a	b	c	d	e		a	b	c	d	e
A	1	0	1	0	1	A	100	100	100	100	67
B	0	1	0	1	0	B	100	100	100	75	50
C	1	0	0	0	0	C	67	67	100	67	100

## Simple solution for Ordinal data

- Let a variable with 3 categories 1, 2, 3  $\Rightarrow (1, 0, 0), (1, 1, 0), (1, 1, 1)$
- Clustering of binary data

## The criterion

$$W(\mathbf{c}, \boldsymbol{\mu}) = \sum_{i=1}^n \sum_{k=1}^g c_{ik}^{\gamma} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2$$

where  $c_{ik} \in [0, 1]$ ,  $\sum_k c_{ik} = 1$ ,  $\sum_i c_{ik} > 0$  and  $\gamma > 1$

## Algorithm

- Randomly select  $K$  objects of  $\Omega$  which form the  $K$  first cluster means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ .
- While not convergence
  - assign each object of  $\Omega$  to the cluster with the nearest cluster mean

$$c_{ik} = \left( \sum_{\ell} \left[ \frac{\|\mathbf{x}_i - \boldsymbol{\mu}_k\|}{\|\mathbf{x}_i - \boldsymbol{\mu}_{\ell}\|} \right]^{\frac{2}{\gamma-1}} \right)^{-1}$$

- The cluster means computed become the new cluster means

$$\boldsymbol{\mu}_k = \frac{\sum_i (c_{ik})^{\gamma} \mathbf{x}_i}{\sum_i (c_{ik})^{\gamma}}$$

# Outline

## 1 Introduction

- Cluster Analysis

## 2 $k$ -means algorithm

- $k$ -means: a family of methods
- Principal points to be retained

## 3 Spherical $k$ -means

## 4 $k$ -means for contingency table

## 5 $k$ -means for categorical data

## 6 Miscellaneous clustering methods

- Sequential methods
- Self-Organizing-Mapping
- LSA, NMF and Clustering
- Dynamic clusters method
- Fuzzy  $k$ -means

## 7 Conclusion

# Conclusion

## Advantages

- Simple and efficient method
- Give readable results
- Complementary to PCA, CA, MCA, MFA etc.
- Extension to contingency tables or categorical data from the principal components\*
- Fuzzy variants of  $k$ -means (see the finite mixture model)
- Methods available in Statistic and data mining Softwares (See R, SAS, etc.)
- Other criteria depending on the nature of data (Document clustering, microarray etc.)

## Disadvantages

- Depend on the shape of clusters
- They require the number of clusters

## Course 3

- Mixture model