

# Unsupervised Classification (Clustering)

Introduction to Machine Learning

Mathilde Mougeot

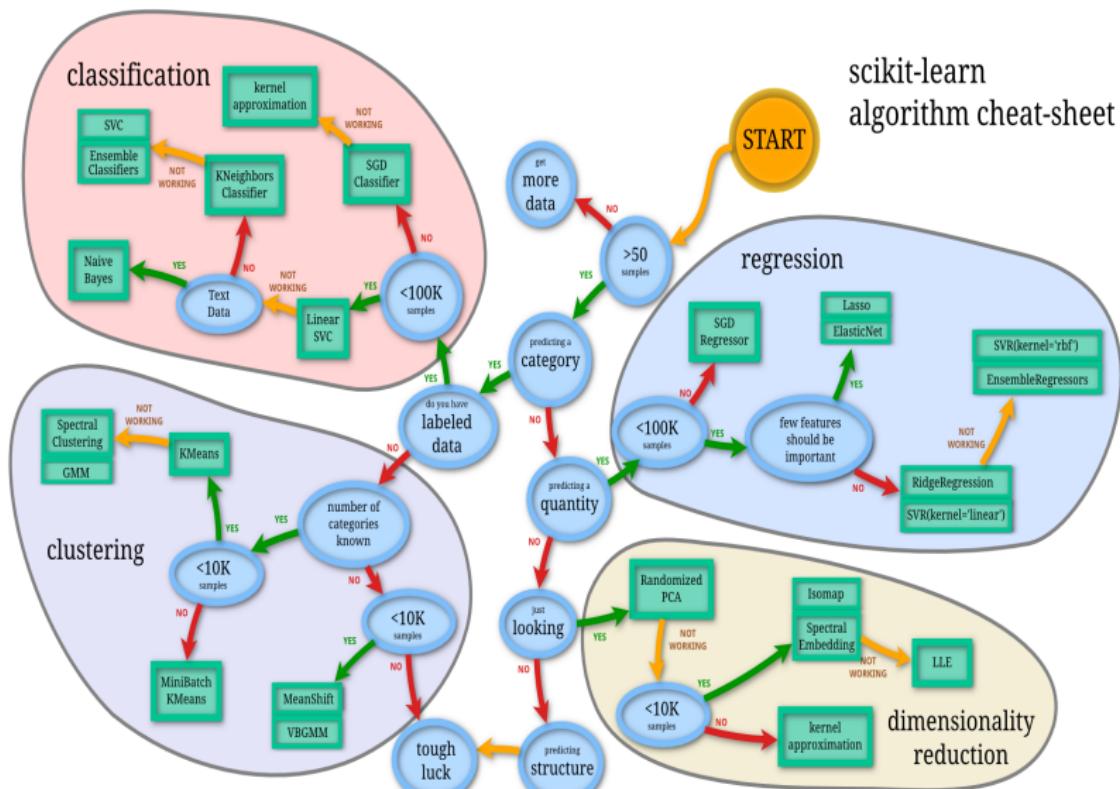
ENSIIE

2020

# At this stage

- Outline
  - Observations  $(X_1, \dots, X_n)$  in  $\mathbb{R}^d$ ,  $d$  may be large
  - Target  $(Y_1, \dots, Y_n)$  labels (classification) or continuous (regression)
  - Goal : Predictive modeling : understand the link between  $X \mapsto Y$ , reduce  $d$ .
- Methods
  - ① **Classification** models (Parametric/ Non Parametric)  
Linear Quadratic Discriminant Analysis, Logistic, KNN CART, Bagging, Random Forest
  - ② **Regression** models (Parametric/ Non Parametric)...  
Linear models, Linear models with penalization : LASSO, Ridge, KNN, CaRt, Bagging, Random Forest...
  - ③ **Clustering** (unsupervised Classification)

# scikit-learn algorithm cheat-sheet



Back

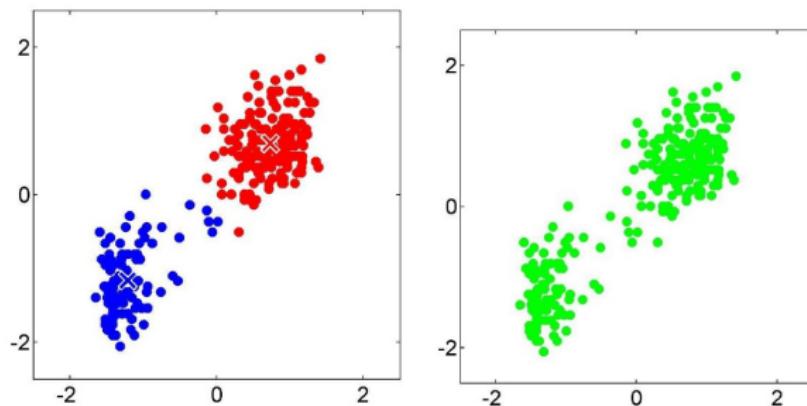
scikit  
learn

# Outline

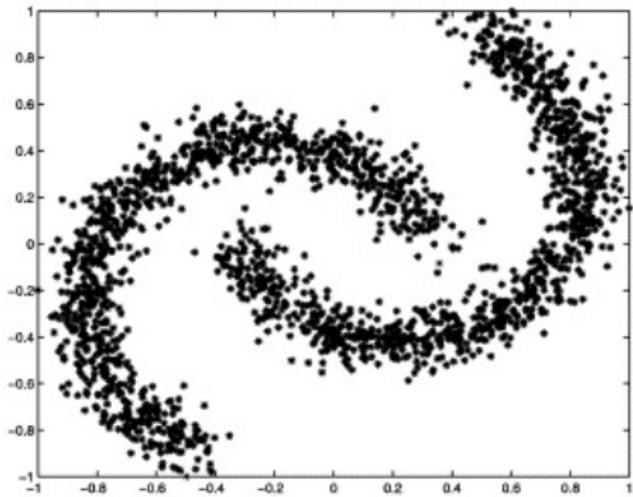
- ① Introduction / motivations
- ② Distance-based clustering - Notations
- ③ Model-based clustering
- ④ Graph-based clustering
- ⑤ Hierarchical clustering
- ⑥ Centroid-based clustering

# Motivations

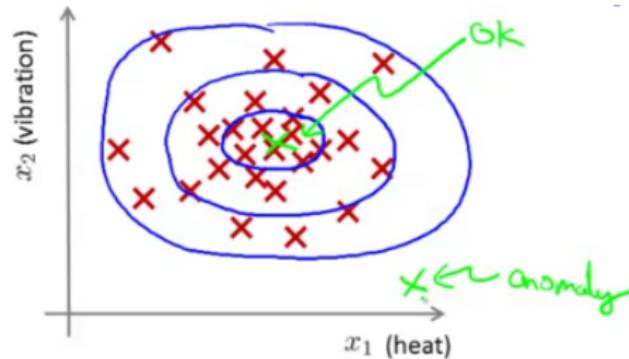
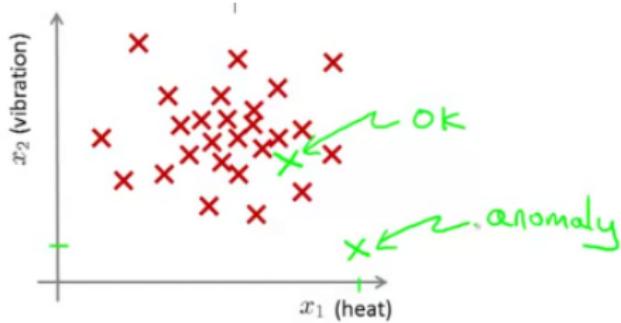
# Unsupervised data (1) - Clustering



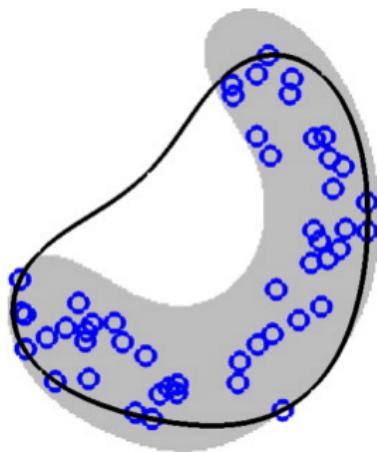
# Clustering can be difficult !



## Unsupervised data (2) - Anomaly/Mode detection



## Unsupervised data (3) - Novelty detection



# Distance-based clustering

-

## Notations

# Clustering input : distance matrix

- Data matrix
  - Individual index  $i \in \{1, \dots, n\}$
  - Feature index  $p \in \{1, \dots, d\}$
  - Measurements  $x_{ip}$
- Distance matrix
  - $p$ -th feature distance between individuals  $i$  and  $j$  =  $d_p(x_i, x_j)$
  - Distance between individuals  $i$  and  $j$  :

$$D(x_i, x_j) = \sum_{p=1}^d w_p d_p(x_i, x_j)$$

where  $w_p$   $p$ -th feature importance,  $w_p > 0$

# Examples of distances

- Quantitative features
  - Squared distance or absolute difference
  - 1-correlation
- Discrete ordinal variables
  - Equidistance encoding
- Categorical variables
  - Zero-one distance
- What if missing values ?

# Cluster dispersion functions

- Encoder function  $C : \{1, \dots, n\} \mapsto \{1, \dots, K\}$  (point to cluster)
- Within-cluster dispersion

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} D(x_i, x_j)$$

- Between-cluster dispersion

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j) \neq k} D(x_i, x_j)$$

- Total dispersion :  $T = W(C) + B(C)$

## Clustering method #1 - Brute force

- Combinatorial assignment
- Number of possibilities for assigning  $n$  points to  $K$  clusters

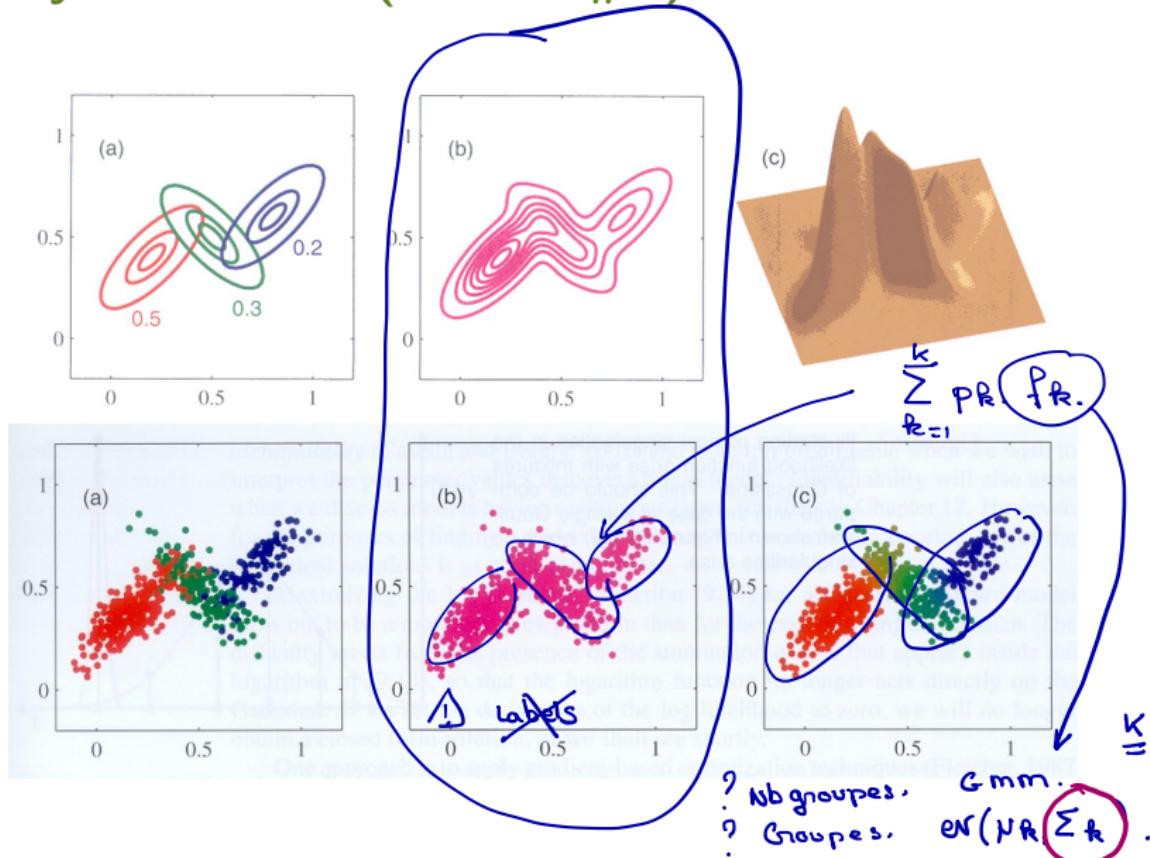
$$S(n, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^n$$

Example :  $S(19, 4) \simeq 10^{10}$

- Question : Limited search vs. approximate solution

# Parametric approach

# Density estimation (Course #1)



## Reminder on gaussian mixture models

- Random vector  $X$  on  $\mathbb{R}^d$  with  $K$  components
- Gaussian densities  $f_k$ ,  $k = 1, \dots, K$ ,
- Component parameters  $(\mu_k, \Sigma_k)$ ,
- Mixture parameter  $p = (p_1, \dots, p_K)$  in the simplex
- Distribution of  $X$  = Mixture density

$$f_X(x) = \sum_{k=1}^K p_k f_k(x) \quad , \quad \forall x \in \mathbb{R}^d$$

- For estimation, use EM algorithm...

## Complexity of the problem depends on...

- The dimension  $d$
- The number of clusters  $K$
- The number of samples  $n$
- The smallest mixture coefficient " $\min_j p_j$ "
- How separated the clusters are...

# High dimension

# Distance between observations

-	$X^1$	$X^2$	...	$X^j$	...	$X^d$
1	$x_{11}$		...	$x_{1j}$		$x_{1d}$
2						
...						
→ $i$	$x_{i1}$		...	$x_{ij}$		$x_{id}$
...						
$n$	$x_{n1}$		...	$x_{nj}$		$x_{nd}$

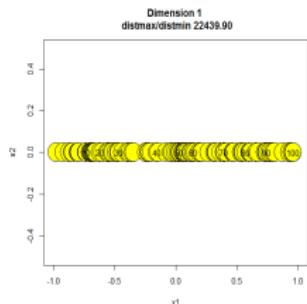
- For two observations  $(x_i, x_k)$ ,  $x_i \in \mathbb{R}^d$ ,  $x_k \in \mathbb{R}^d$
- Euclidian distance  $\ell_2$  between two observations

$$\|x_i - x_k\|_{\ell_2} = \sqrt{\sum_{j=1}^d (x_i(j) - x_k(j))^2}$$

# Dimension curse

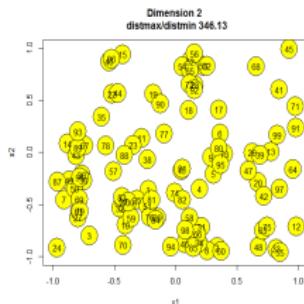
- Evaluation of the distance between two observations in dimension  $d$
- Illustrations :  $n = 100$  observations uniformly distributed, 1, 2, 3, ...
- Indicator :  $\frac{\max_{i \neq j} \|x_i - x_j\|_{\ell_2}}{\min_{i \neq j} \|x_i - x_j\|_{\ell_2}}$

$$d = 1$$



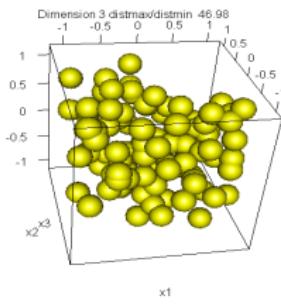
22 435

$$d = 2$$



346

$$d = 3$$

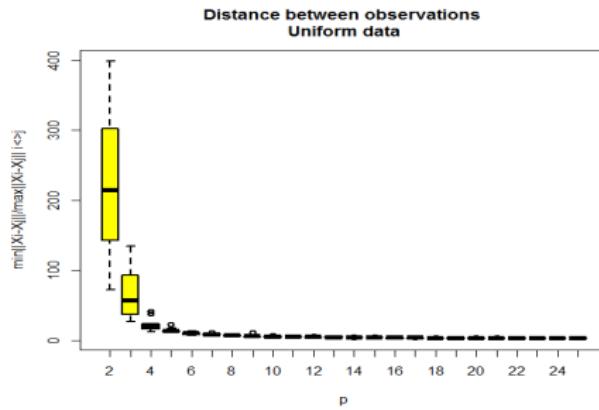


47

## Dimension curse

Ratio study  $\frac{\max_{i \neq j} \|x_i - x_j\|_{\ell_2}}{\min_{i \neq j} \|x_i - x_j\|_{\ell_2}}$  function of the dimension  $d$

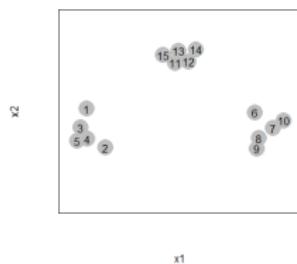
Illustration :  $n = 100$  observations uniformly distributed ( $K = 100$  repetitions)



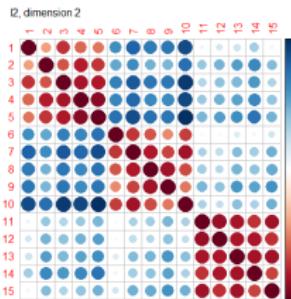
- The value of the ratio tends to  $\sim 1$  when  $d$  increases.
- The euclidian distance loses its discrimination ability in high dimension
- Serious problem especially for segmentation tasks...

# Data segmentation ( $d=2$ )

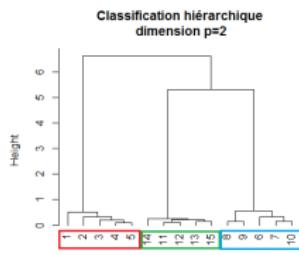
Observations



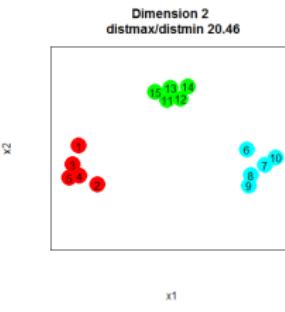
distance matrix



HAC

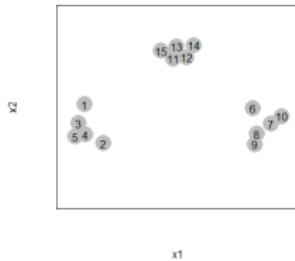


3 classes Clustering

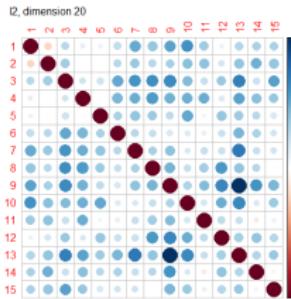


# Data segmentation ( $d=20$ )

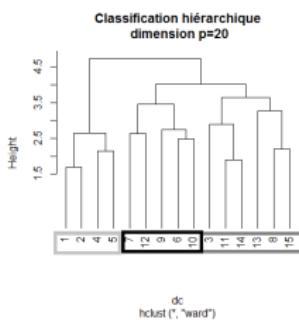
data are embedded in a high dimensional space  $d = 20 = 2 + 18$   
Observations



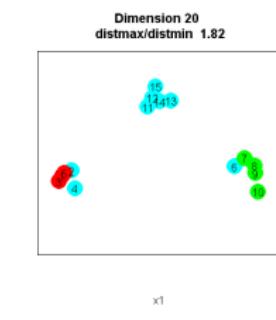
distance matrix



HAC



3 classe Clustering



MAL

## Dimension reduction

Find good representations of the data initially coded in large dimensions

- **Features** : a small number of discriminant features based on data expertise or automatic extraction.
- **Compress Sensing** : sparse representation ( $S$ ) of  $x$  based on a linear combination of  $p$  vectors.
- **Manifold estimation** :  $x$  is represented in a low-dimensional space using the Laplacian eigenvectors on the variety, estimated from a graph of neighborhoods using the examples

→ Mathematical tools at the interface of harmonic analysis, geometry, probability and statistics.

# Model based Clustering

# Model-based clustering

Set of observations  $\{x_i, x_i \in \mathbb{R}^d, 1 \leq i \leq N\}$

**Assumptions :** Mixture of  $K$  gaussian :  $f_X(x) = \sum_{k=1}^K \pi_k f_k(x)$   
 $\mu_k$  (means),  $\Sigma_k$  (covariances),  $\pi_k$  (mixing coefficients)  
 $1 \leq k \leq K$

**Find ? :**  $\mu_k, \Sigma_k, \pi_k$

using the EM Algorithm :

- ① Initialization
- ② E Step : Expectation Step
- ③ M Step : Maximization Step
- ④ LogLikelihood computation

# EM for gaussian mixture (1/4)

## ① Initialization :

$\mu_k$  (means),  $\Sigma_k$  (covariances),  $\pi_k$  (mixing coefficients)  
Compute Log Likelihood

$$\ln p(X|\mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\}$$

## ② E Step :

## ③ M Step :

## ④ Evaluate the log likelihood :

## EM for gaussian mixture (2/4)

- ① **Initialization** :  $\mu_k$ ,  $\Sigma_k$ ,  $\pi_k$ , Compute Log Likelihood
- ② **E Step** : Evaluate the responsibilities using current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}$$

- ③ **M Step** :
- ④ **Evaluate the log likelihood** :

# EM for gaussian mixture (3/4)

① **Initialization** :  $\mu_k$ ,  $\Sigma_k$ ,  $\pi_k$ , Compute Log Likelihood

② **E Step** :  $\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}$

③ **M Step** : Re-estimate the parameters using the current responsibilities :

- $\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$
- $\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k^{new})(x_n - \mu_k^{new})^T$
- $\pi_k^{new} = \frac{N_k}{N}$  where  $N_k = \sum_{n=1}^N \gamma(z_{nk})$

④ **Evaluate the log likelihood** :

# EM for gaussian mixture (4/4)

① **Initialization** :  $\mu_k$ ,  $\Sigma_k$ ,  $\pi_k$ , Compute Log Likelihood

② **E Step** :  $\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}$

③ **M Step** : Re-estimate the parameters using the current responsibilities :

- $\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$
- $\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k^{new})(x_n - \mu_k^{new})^T$
- $\pi_k^{new} = \frac{N_k}{N}$  where  $N_k = \sum_{n=1}^N \gamma(z_{nk})$

④ **Evaluate the log likelihood** :

$$\ln p(X | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\}$$

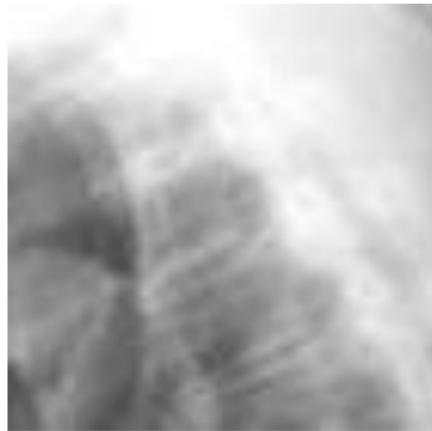
REPEAT 2,3,4 UNTIL CONVERGENCE

## Model-based clustering

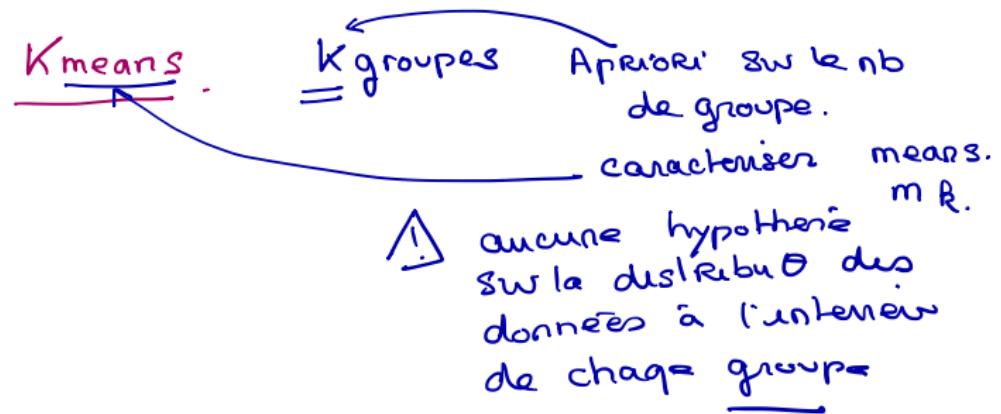
- The mixture density assumption (gaussian) should be valid
- How to chose the number of clusters ? (value of  $K$  )  
penalization of the likelihood

# Illustration

Image segmentation based on model-based clustering



# Centroid-based clustering



# The celebrated K-means (1)

$$\mathcal{D}_n = \{x_i, \quad x_i \in \mathbb{R}^d, 1 \leq i \leq n\}$$

- Use for : Quantitative features
- Squared Euclidean distance
- Barycenter cluster  $k$  with  $n_k$  elements

$$\bar{x}_k = n_k^{-1} \sum_{C(i)=k} x_i$$

$n_k$  : nb d'observations  
 $\in C_k$ .

- Note that : for any subset  $I \subset \{1, \dots, n\}$  of individuals

$$\bar{x}_I = \arg \min_m \sum_{i \in I} \|x_i - m\|^2$$

## The celebrated $K$ -means (2)

- Optimization criterion

1. . . ns nble des groupes  $K$ ,

$$W(C) = \sum_{k=1}^K \sum_{C(i)=k} \|x_i - \bar{x}_k\|^2$$

- Solution

$$C^* = \arg \min_C W(C) = \arg \min_{C, m_1, \dots, m_K} \sum_{k=1}^K \sum_{C(i)=k} \|x_i - m_k\|^2 + \lambda \overset{?}{\circ}(k)$$

$m_1^*, C_1^*$   
 $m_2^*, C_2^*$   
 $\vdots$   
 $m_K^*, C_K^*$

# Clustering method #3a - K-means algorithm

on suppose  $K$  connu.

Parameter : encoder range  $K$

**Initialization** : initial encoder  $C^{(0)}$ , and centers  $m_k^{(0)}$

$c_1^{(0)}, c_2^{(0)} \dots c_K^{(0)}$ .  
 $m_1^{(0)}, m_2^{(0)} \dots m_K^{(0)}$ .

**Step 1** : Fix encoder  $C$ , compute the centers

$$\bullet \quad m_k = \frac{1}{n_k} \sum_{C(i)=\underline{k}} x_i$$

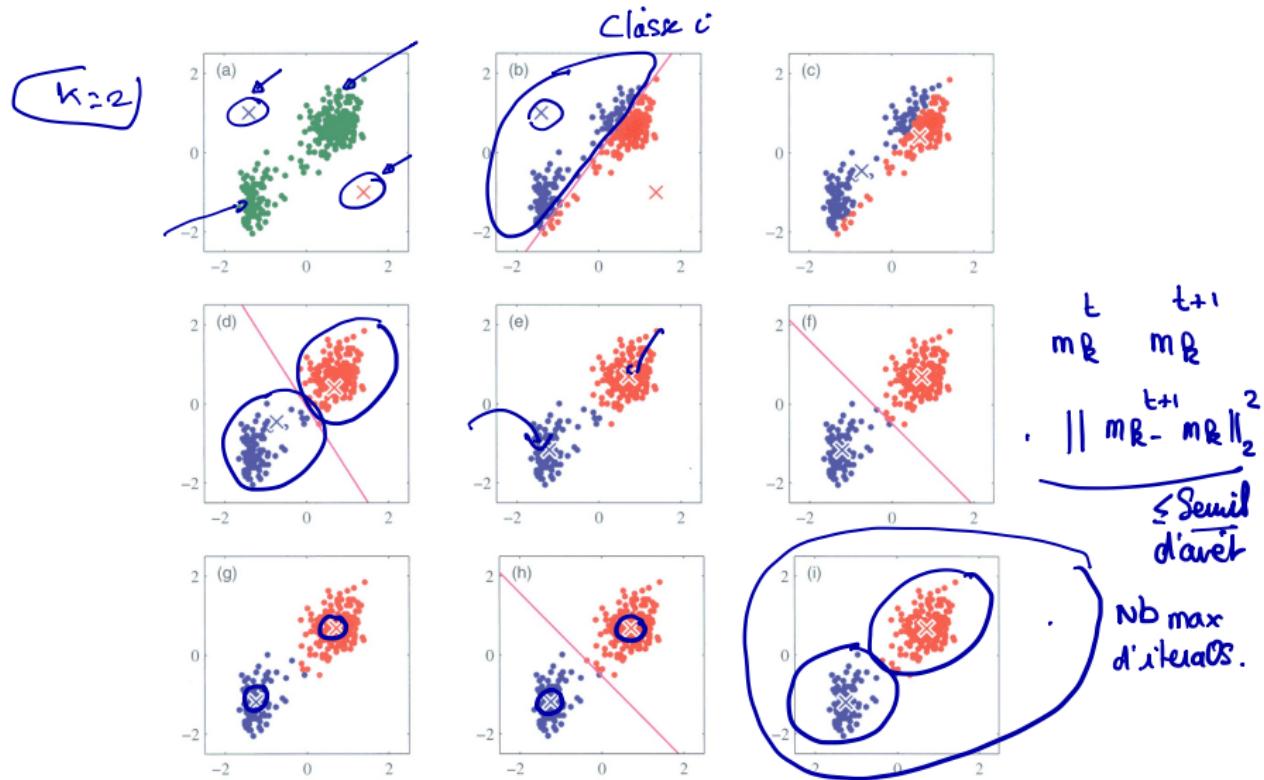
**Step 2** : Fix the centers  $m_1, \dots, m_K$ , assign with new encoder

(\*)  $\underline{C(i)} = \arg \min_{1 \leq k \leq K} \|x_i - \underline{m_k}\|^2$

**Iteration** : repeat Steps 1 & 2 until  $C$  does not change anymore.

==. stabilité des groupes  
. > Nb d'itérations fixées.

# K-means algorithm - How it works



## A variation : K-medoids

point "le plus profond."  
ex: mediane en dim 1.

- Use for : any type of features
- Arbitrary distance
- Centers  $\{m_1, \dots, m_K\}$  belong to the data set  $\{x_1, \dots, x_n\}$

## Clustering method #3b - $K$ -medoids algorithm

**Parameter** : encoder range  $K$

**Initialization** : initial encoder  $C^{(0)}$ , and centers  $m_k^{(0)}$

**Step 1** : Fix encoder  $C$ , compute the centers  $m_k = x_{i_k^*}$  with

$$i_k^* = \arg \min_{C(i)=k} \sum_{C(j)=k} D(x_i, x_j)$$

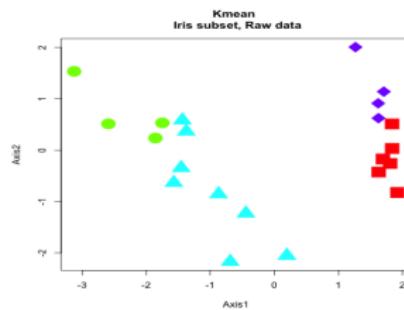
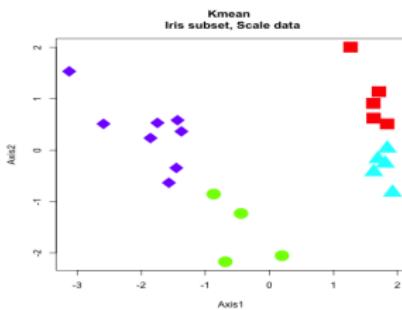
**Step 2** : Fix the centers  $m_1, \dots, m_K$ , assign with new encoder

$$C(i) = \arg \min_{1 \leq k \leq K} D(x_i, m_k)$$

**Iteration** : repeat Steps 1 & 2 until  $C$  does not change anymore.

# Kmeans Clustering. Illustration

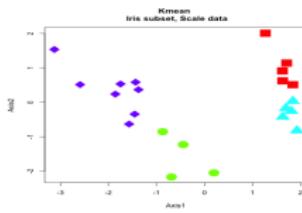
Use of PCA to project the observations, and to represent the clusters.



Subset (#30) of Iris Data Set  
Impact of Scaling data  
(left :scaled vs right :raw)

## Kmeans Clustering. R instructions.

```
# tab : dataframe  
pca=dudi.pca(tab,scannf=FALSE,nf=2);  
K=4;  
mycol=rainbow(K); mypch=c(1,3,4,8)  
#Scale data  
res=kmeans(tab,centers=4);  
plot(pca$li,col=mycol[res$cluster],pch=mypch[res$cluster])
```

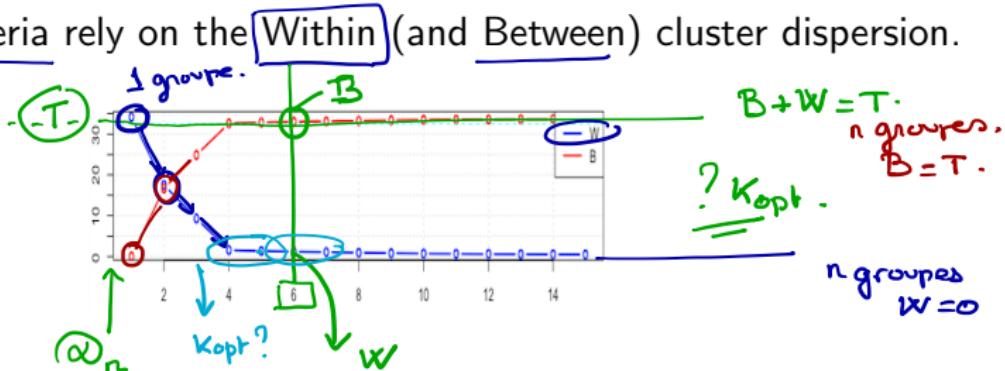


# Calibrating the number of clusters

$K^?$

# Calibrating the number of clusters

Most of the criteria rely on the Within (and Between) cluster dispersion.



- Encoder function  $C : \{1, \dots, n\} \mapsto \{1, \dots, K\}$  (point to cluster)

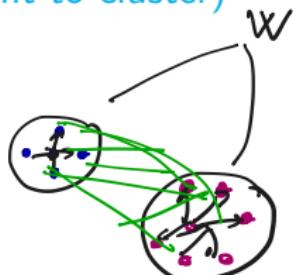
Total dispersion :  $T = W(C) + B(C)$

- Within-cluster dispersion :

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} D(x_i, x_j)$$

- Between-cluster dispersion :

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j) \neq k} D(x_i, x_j)$$



. B/w  
 . stability  
 B  
 w

## Calibrating the number of clusters (1)

- Calinski & Harabasz (1974)

$$\text{maximize } F_{CH}(K) = \frac{B(C_K)/(K-1)}{W(C_K)/(n-K)}, \quad \forall K > 1.$$

- Hartigan (1975)

take smallest  $K \geq 1$  such that  $F_H(K) \leq 10$ ,

where

$$F_H(K) = \left( \frac{W(C_K)}{W(C_{K+1})} - 1 \right) / (n - K - 1).$$

Evoluo de la variancia w.

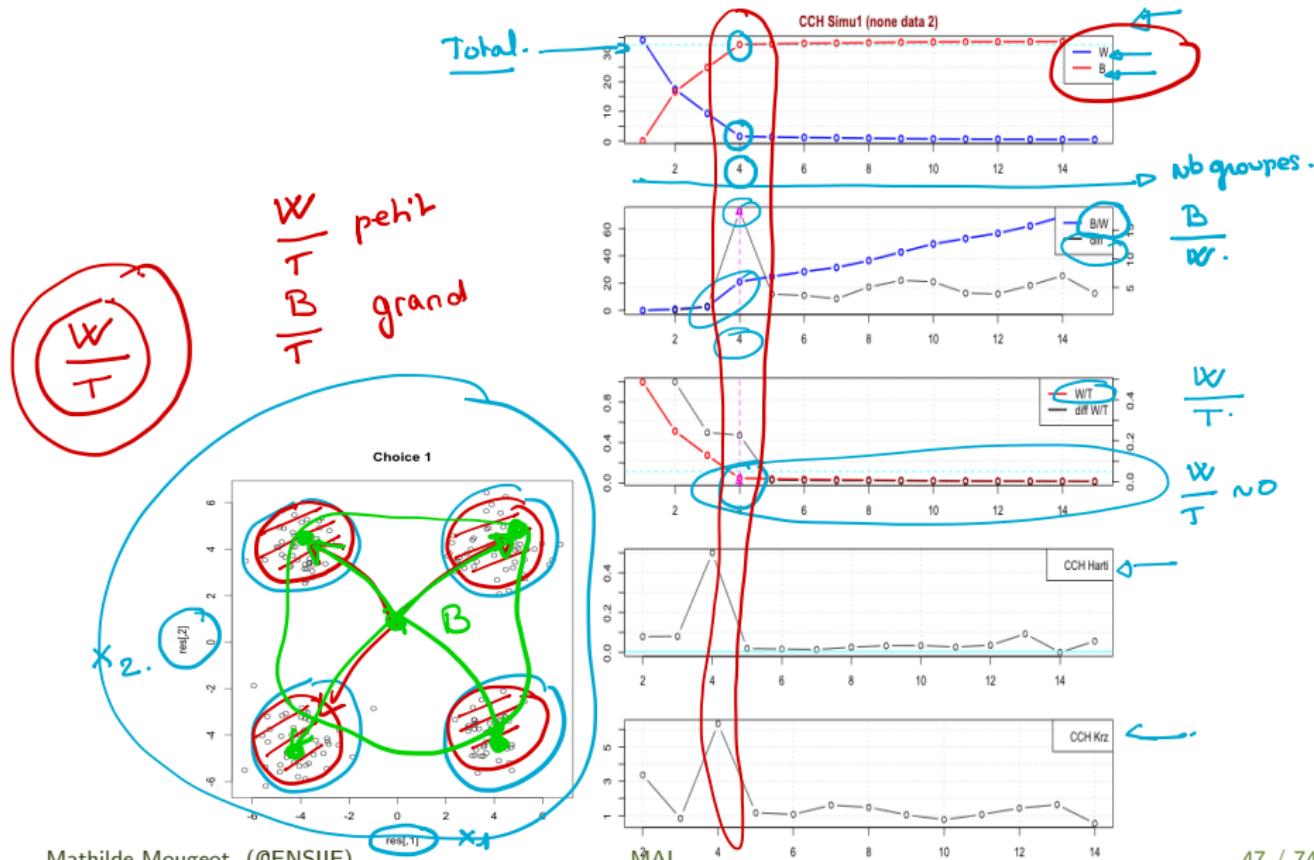
## Calibrating the number of clusters (2)

- Krzanowski & Lai (1985)

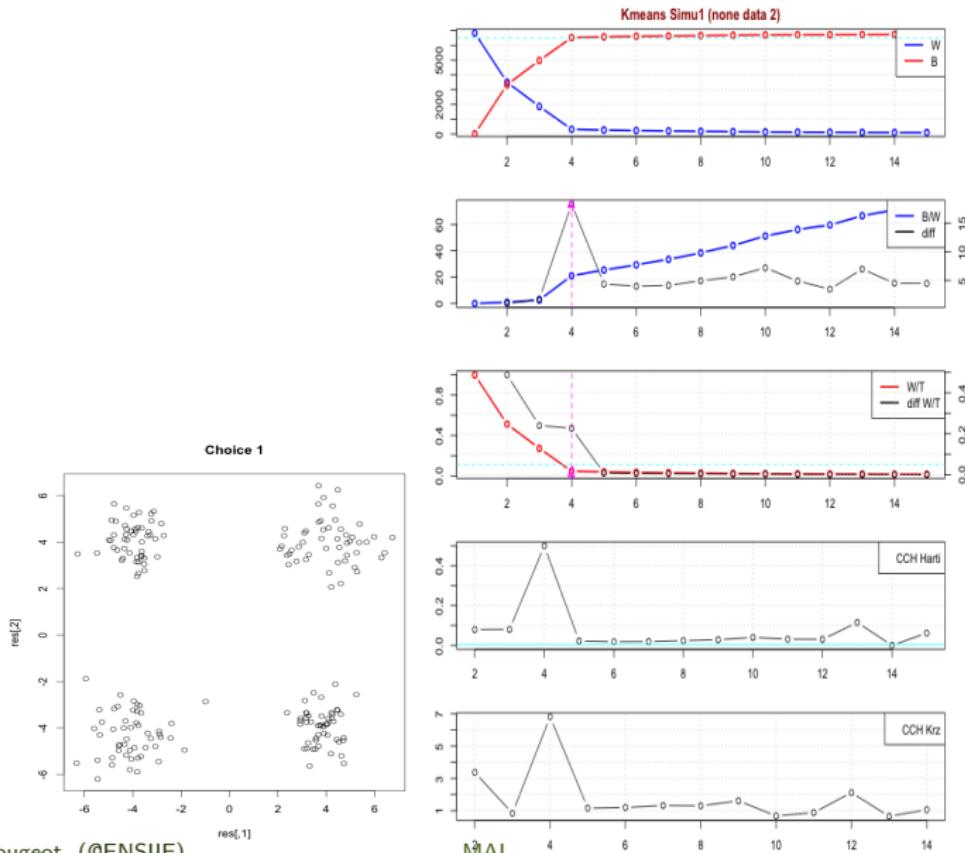
$$\text{maximize } F_{KL}(K) = \left| \frac{\Delta(K)}{\Delta(K+1)} \right|$$

where  $\Delta(K) = (K-1)^{2/d} W(C_{K-1}) - K^{2/d} W(C_K)$  and  $d$  is the dimension of input data.

# Illustration. Calibrating the number of clusters



# Illustration. Calibrating the number of clusters



## Calibrating the number of clusters (3)

- Rousseeuw (1987) - Silhouette statistic indicate how mixed.

$$F_S(K) = \sum_{i=1}^n \left( \frac{b(i) - a(i)}{\max(a(i), b(i))} \right),$$

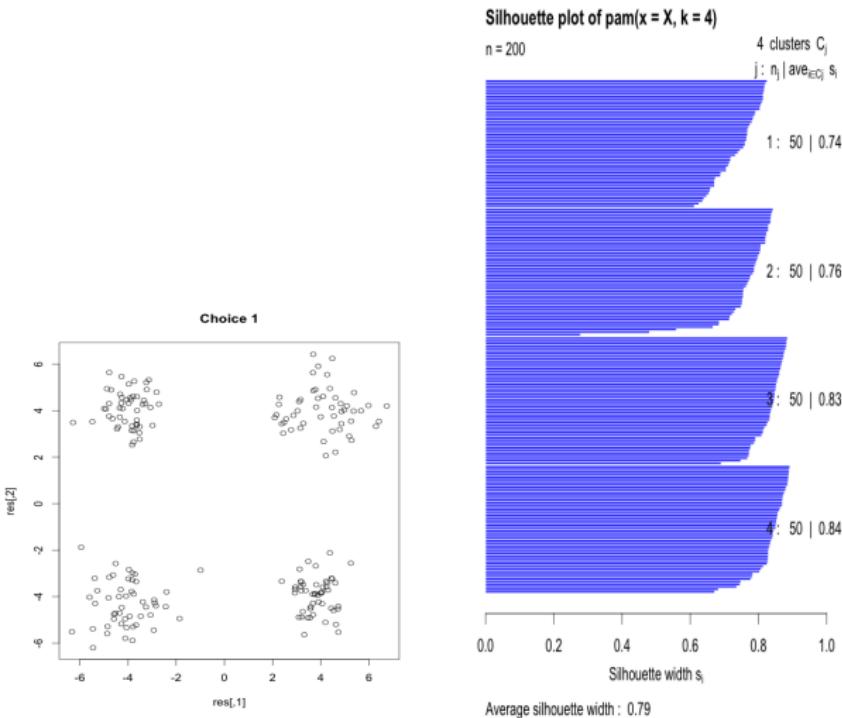
where for  $i \in C_k$

$$a(i) = \frac{1}{n_k} \sum_{C(j)=k} \|x_j - x_i\|_2^2,$$

and, if  $\ell = \ell(i)$  is the next nearest cluster of the point  $x_i$  :

$$b(i) = \frac{1}{n_\ell} \sum_{C(j)=\ell} \|x_j - x_i\|_2^2.$$

# Calibrating the number of clusters. Illustration.



# Calibrating the number of clusters. Illustration.

Silhouette plot of  $\text{pam}(x = X, k = 4)$

$n = 200$

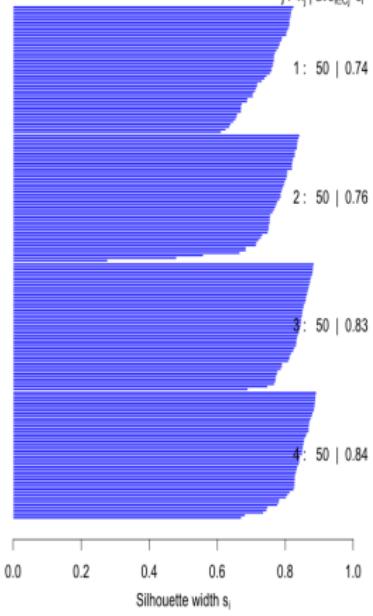
4 clusters  $C_j$   
 $j : n_j | \text{ave}_{EC_j} s_i$

1: 50 | 0.74

2: 50 | 0.76

3: 50 | 0.83

4: 50 | 0.84



Average silhouette width : 0.79

Silhouette plot of  $\text{pam}(x = X, k = 3)$

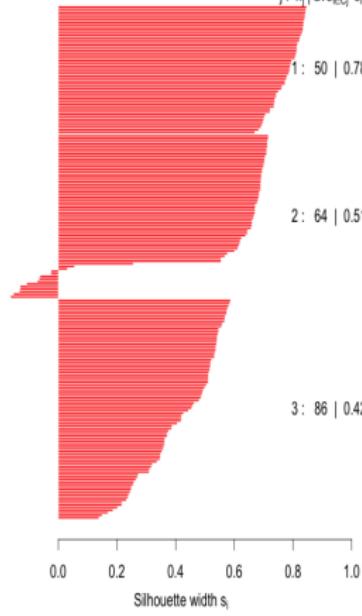
$n = 200$

3 clusters  $C_j$   
 $j : n_j | \text{ave}_{EC_j} s_i$

1: 50 | 0.78

2: 64 | 0.51

3: 86 | 0.42



Average silhouette width : 0.54

Silhouette plot of  $\text{pam}(x = X, k = 5)$

$n = 200$

5 clusters  $C_j$   
 $j : n_j | \text{ave}_{EC_j} s_i$

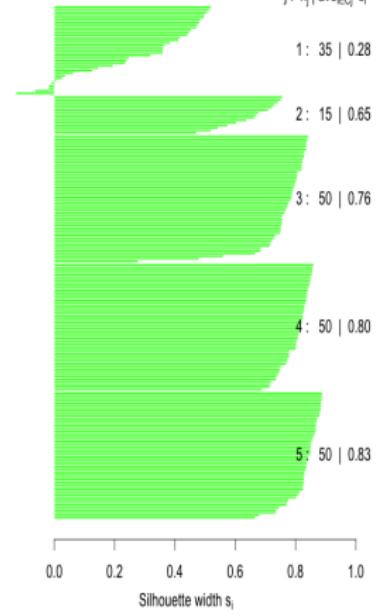
1: 35 | 0.28

2: 15 | 0.65

3: 50 | 0.76

4: 50 | 0.80

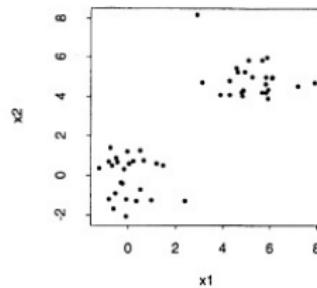
5: 50 | 0.83



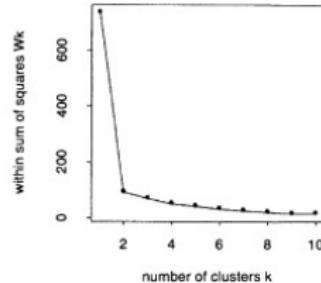
Average silhouette width : 0.7

# Calibrating the number of clusters (4)

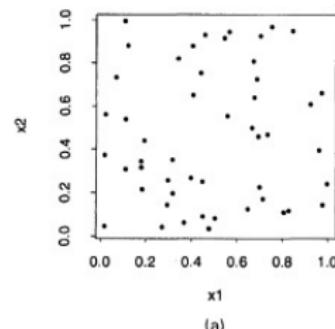
- Tibshirani, Walther, Hastie (2001) - Gap statistic



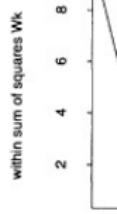
(a)



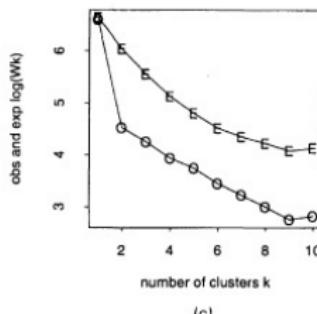
(b)



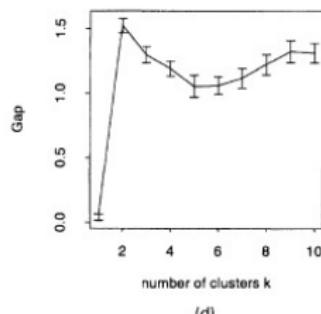
(a)



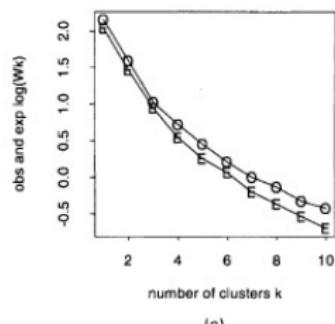
within sum of squares  $W_k$



(c)



(d)



(e)

# Theoretical analysis

- An information-theoretic formulation
- References on Vector Quantization
- Work of Pollard (1981, 1982), but also Linder (2001)
- Proofs of strong consistency of  $K$ -means clustering

# K-means clustering vs EM Algorithm

Algo ↗

Algo mélange Gaussien.

$G^{mm}$

- Distortion measure :

$$J = \sum_{i=1}^n \sum_{k=1}^k r_{ik} \|x_i - \mu_k\|^2$$

Matrice de Responsabilité  
 $\gamma [n, k]$   
 $\uparrow \quad \uparrow$

$x_i \in \mathcal{R}^d$   
 $r_{ik}$   
 $\mu_k$

$$r_{ik} = 1 \text{ if } k = \arg \min_j \|x_i - \mu_j\|^2 \quad (=0 \text{ otherwise})$$

- K-means algorithm :

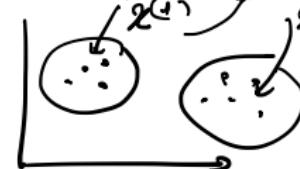
- ① E-Step :  $r_{ik}$  computation
- ② M-step :  $\mu_k$  computation

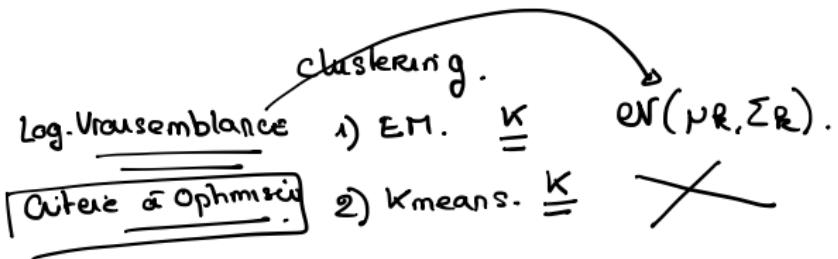
REPEAT 1, 2 UNTIL convergence

$$r_{ik} \quad 1 \leq i \leq n.$$

$$r_{ik} = 1 \text{ si } x_i \in C_k.$$

$$x_i \in \mathcal{R}^d = \bar{x} \in \mathcal{R}^d.$$



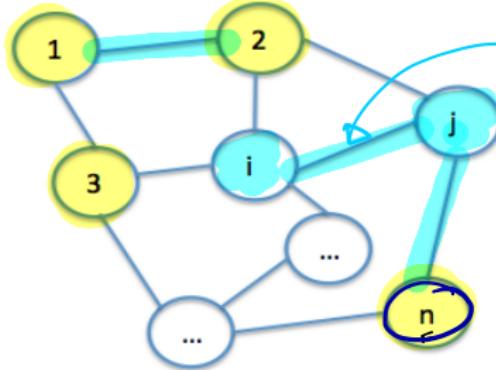


# Graph-based clustering

clustering Spectral

# A Graph

$$W = \begin{bmatrix} & & j \\ i & & \\ & w_{ij} & \end{bmatrix}$$



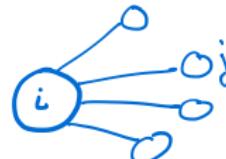
→ clustering sur  $\mathcal{Q}_{n-d}$   
 $\mathcal{Q}_n = \{x \in \mathbb{R}^d\}$   
à partir de  $n$  points

$w_{ij}$  poids de connexion entre 2 noeuds du graphe -

- $w_{ij} > 0$
- $w_{ij} = 0$  pas de lien direct entre  $i$  et  $j$ .

- $G = (V, E)$  is a graph with vertex set  $V = \{v_1, \dots, v_n\}$  and Edges.
- $w_{i,j}$  : positive weight between node  $i$  and node  $j$ .  
If  $w_{ij} = 0$  then vertices  $v_i$  and  $v_j$  are not connected
- Weighted adjacency matrix  $W = (w_{ij})_{1 \leq i,j \leq n}$  with positive coefficients
- Undirected graph means  $W$  symmetric :  $w_{i,j} = w_{j,i}$ .

# Graph definitions

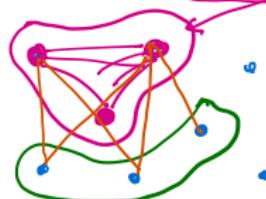


- The degree of a vertex (node) of index  $i$  defined by :  $\deg_i = \sum_{j=1}^n w_{ij}$
- The Degree matrix is a diagonal matrix defined by  
 $D = \text{diag}(\deg_1, \dots, \deg_n)$

$$D = \begin{bmatrix} & & & \\ & d_{11} & & \\ & & d_{22} & \\ & & & d_{nn} \end{bmatrix}$$

di<sub>ij</sub> = 0 if i ≠ j

- For the subset  $A \subset \{1, \dots, n\}$  :



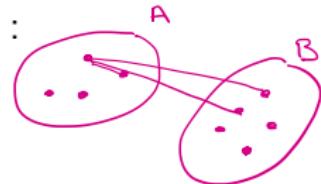
$$\begin{aligned} |A| &= \text{cardinality of } A \quad \# \text{ noeuds de } A. \\ \text{vol}(A) &= \sum_{i \in A} \deg_i \end{aligned}$$

- Let  $A$  and  $B$  two disjoint subsets of nodes of  $\{1, \dots, n\}$

The MinCut distance between  $A$  and  $B$  is defined by :

$$W(A, B) = \sum_{i \in A, j \in B} w_{ij}$$

```
graph LR; subgraph A; ...; end; subgraph B; ...; end; A --- B
```



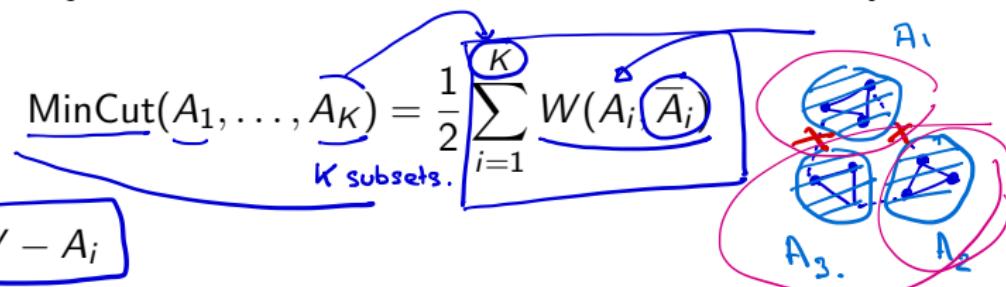
# Graph cut formulation

In order to understand the underlaying structure of the graph, we are interested in discovering subsets of nodes characterized by edges within group with high weights and edges between groups with low weights and edges.

- Considering K disjoint subsets, the MinCut criterion is defined by :

$$\text{MinCut}(A_1, \dots, A_K) = \frac{1}{2} \sum_{i=1}^K W(A_i; \bar{A}_i)$$

where  $\bar{A}_i = V - A_i$



- Drawback : Often leads to a cluster such that  $|A_1| = 1$  if  $K = 2$ .

# Alternative criteria to MinCut

- To guarantee large enough clusters, other normalized criteria are introduced :

Critère :

MinCut

$$\rightarrow \text{RatioCut}(A_1, \dots, A_K) = \frac{1}{2} \sum_{i=1}^K \frac{W(A_i, \bar{A}_i)}{|A_i|}$$

# noeuds.

$$\rightarrow \text{NCut}(A_1, \dots, A_K) = \frac{1}{2} \sum_{i=1}^K \frac{W(A_i, \bar{A}_i)}{\text{vol}(A_i)}$$

Somme des poids des noeuds du groupe.

- Drawback : NP-hard optimization problems

pb d'optimisation non réalisable lorsque  $n$  est grand.

# Spectral clustering

- The Spectral clustering algorithm proposes a relaxation of the RatioCut and Ncut minimizations.
- The computation of the eigenvectors of the Graph Laplacian operator approximate the RatioCut solution.

# Graph Laplacian (1)

Notations :

$W$  is the adjacency matrix  $W = (w_{ij})_{1 \leq i,j \leq n}$  with  $w_{ij} \geq 0$ .

$D$  =  $\text{diag}(\deg_1, \dots, \deg_n)$  is the diagonal degree matrix :  $\deg_i = \sum_{j=1}^n w_{ij}$ .

- Definition - Unnormalized graph Laplacian matrix :

$$L = D - W$$

- Property 1 - For any vector  $f \in \mathbb{R}^n$

$$\begin{aligned} f^T L f &= f^T (D - W) f \\ &= f^T D f - f^T W f \\ &= \sum_{i=1}^n d_i \cdot f_i^2 - \sum_{i,j} w_{ij} f_i f_j \end{aligned}$$

$$f^T L f = \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2$$

- Property 2 -  $L$  symmetric, positive

- Property 3 - Smallest eigenvalue of  $L$  is 0

- Property 4 - Relation between the spectrum of  $L$  and the number of connected components of the graph.

$$\begin{aligned} f^T L f &= \\ &= \sum_{i=1}^n d_i \cdot f_i^2 - \underbrace{\sum_{i,j} w_{ij} f_i f_j}_{f^T W f} \\ &= \sum_{i=1}^n d_i \cdot f_i^2 - \sum_{i,j} f_i f_j w_{ij} \\ &= \frac{1}{2} \left( \sum_{i=1}^n d_i f_i^2 - 2 \sum_{i,j} f_i f_j w_{ij} + \sum_{i=1}^n d_i f_i^2 \right) \\ &= \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2 \end{aligned}$$

(values prop. -  
vect. prop. - de  $L$ .)

## Building a similarity graph

$$\mathcal{Q}_n = \{x_i \mid 1 \leq i \leq n \quad x_i \in \mathbb{R}^d\}$$

calcul  $s(x_i, x_j)$ .  $\|x_i - x_j\|_2^2$  power.

In practice :

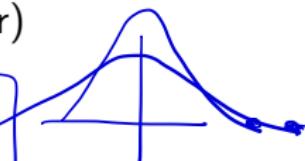
matrice de  
similarité.

Input : Similarities  $s(x_i, x_j)$  or distances  $D(x_i, x_j), \forall i, j$

- Following options are often applied on the similarity matrix, and appeared to be very useful based on the (introduction of thresholding or non-linearities)

- Option 1 -  $\epsilon$ -neighborhood graph, ( $\epsilon$  hyper parameter)
- Option 2-  $k$ -nearest neighbor graph, ( $k$  hyper parameter)
- Option 3 - Fully connected graph, ( $\sigma$  hyper parameter)

$$s(x_i, x_j) = \exp(-\|x_i - x_j\|^2 / (2\sigma^2))$$



# Clustering method - Unnormalized Spectral Clustering Algorithm

**Input** : Similarity matrix  $S$ . The number  $K$  of clusters is given (as Kmeans)

Qn. (opt1, opt2, opt3)

**Preprocessing** : a similarity graph with adjacency matrix  $W$  and Compute the unnormalized Laplacian  $L$

**Solve eigenvalue problem** : compute the first  $K$  eigenvectors of

$$L = D - W$$

à partir d'une nouvelle représentation des points.

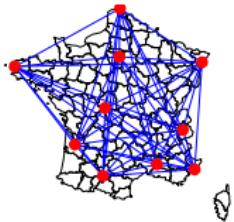
**Clustering in feature space** : Let  $U \in \mathbb{R}^{n \times K}$  be the matrix containing the vectors  $u_1, \dots, u_K$  as columns (eigenvectors with normalisation)

$\begin{bmatrix} \vdots & \\ u_1 & \dots & u_K \end{bmatrix}_n$

- For  $i = 1, \dots, n$ , let  $y_i \in \mathbb{R}^K$  be the vector corresponding to the  $i$ -th row of  $U$
- Cluster the points  $(y_i)_{i=1,\dots,n}$  in  $\mathbb{R}^K$  with the  $K$ -means algorithm into clusters encoded by a partition  $A_1, \dots, A_K$

# Spectral clustering

Full connected graph with  $n$  nodes.



$$\mathcal{Q}_n = \{z_i \mid z_i \in \mathbb{R}^d\}$$

$n=10$

$z_i$  : feature vector  
aux others zero.

Weight between two nodes

$(Z_i, Z_j)$  :

$$\frac{-||z_i - z_j||_2^2}{2\mu^2}$$

$\mu$  heat parameter

Normalized Graph Laplacian :

$$L = I - D^{-1/2} W D^{-1/2}$$

$$L' = D^{-1/2} L D^{-1/2} \\ = I -$$

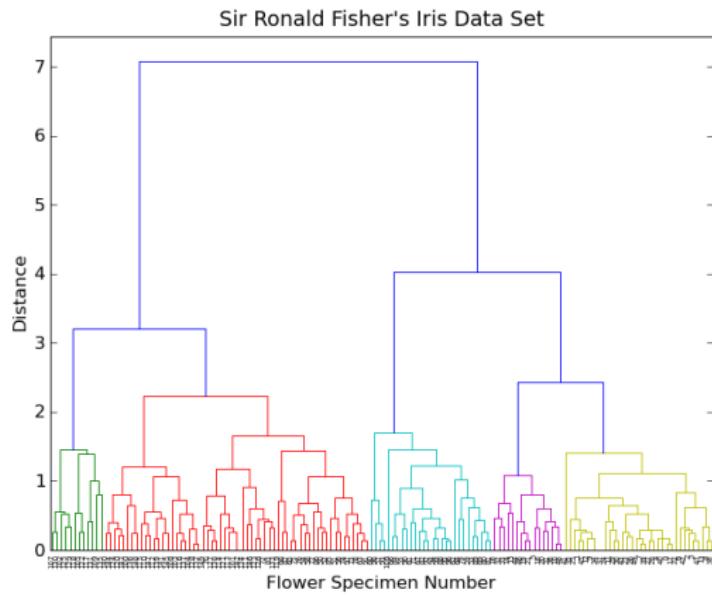
Ng et al. Algorithm (2002) :

**Input :** Fix  $k$  nb . clusters

- 1 Compute the first  $k$  eigenvectors  $u_1, \dots, u_k$  of  $L$  corresponding to the "k" smallest eigenvalues,
- 2 let  $U \in \mathbb{R}^{n \times k}$  be the matrix of column vectors  $u_1, \dots, u_k$
- 3 Form the matrix  $T \in \mathbb{R}^{n \times k}$   
 $t_{i,j} = u_{i,j} / (\sqrt{\sum_k u_{ik}^2})$ .  
Let  $y_i \in \mathbb{R}^k$   $i^{th}$  row of  $T$ .
- 4 Cluster  $\{y_i\}$ ,  $1 \leq i \leq n$  with the **k-means** into clusters  $C_1, \dots, C_k$

# Hierarchical clustering

# Dendrogram



## Clustering method #2a - Bottom-up heuristic

**Input** : number  $K$  of clusters, distance matrix  $D(x_i, x_j)$ , cluster distance  $\Delta$

**Initial step** : find the pair  $(x_i, x_j)$  the minimal element in the distance matrix and form cluster  $A_1 = \{i, j\}$ , the remaining  $x_k$ 's form clusters with one element  $A_2, \dots, A_{n-1}$

**Agglomeration step** : Consider  $A_1, \dots, A_{n-1}$  clusters and find the pair  $(k^*, l^*)$  such that

$$(k^*, l^*) = \arg \min_{k \neq l} \Delta(A_k, A_l)$$

and merge these clusters into  $A = A_{k^*} \cup A_{l^*}$ .

**Stopping criterion** : Iterate until the target number of clusters is reached.

# Linkage distance

- $A, B \subset \{1, \dots, n\}$
- Single linkage

$$\Delta(A, B) = \min_{i \in A, j \in B} \{D(x_i, x_j)\}$$

- Complete linkage

$$\Delta(A, B) = \max_{i \in A, j \in B} \{D(x_i, x_j)\}$$

- Centroid linkage

$$\Delta(A, B) = D(\bar{x}_A, \bar{x}_B)$$

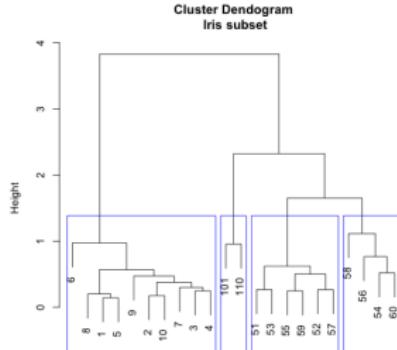
- Average linkage

$$\Delta(A, B) = \frac{1}{|A| \cdot |B|} \sum_{i \in A} \sum_{j \in B} D(x_i, x_j)$$

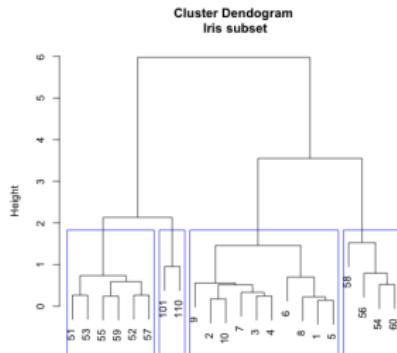
- Ward linkage

$$\Delta(A, B) = \frac{|A| + |B|}{|A| \cdot |B|} D(\bar{x}_A, \bar{x}_B)$$

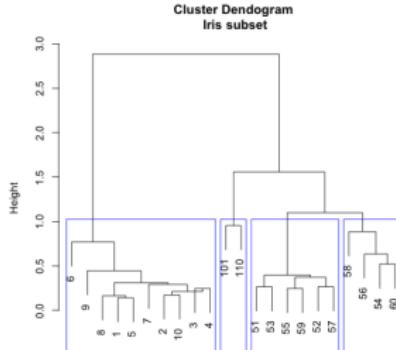
# Hierachical Clustering. Impact of Linkage. Illustration



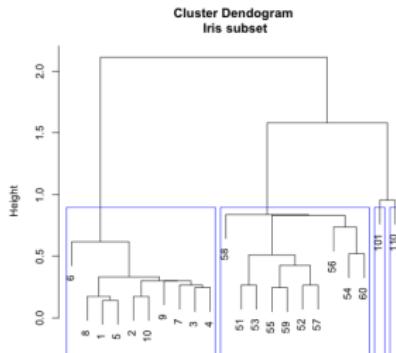
$D$   
hclust(\*, "average")



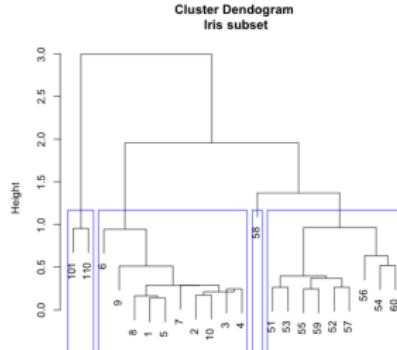
$D$   
hclust(\*, "complete")



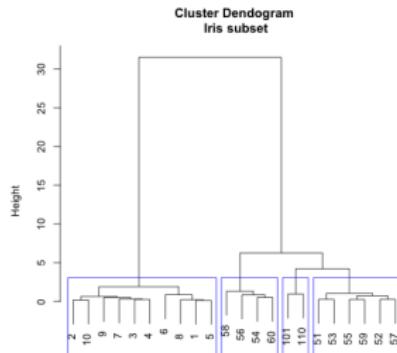
$D$   
hclust(\*, "centroid")



$D$   
hclust(\*, "single")

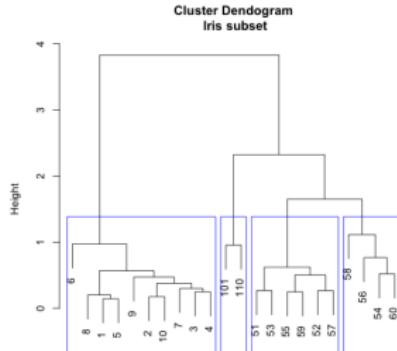


$D$   
hclust(\*, "median")

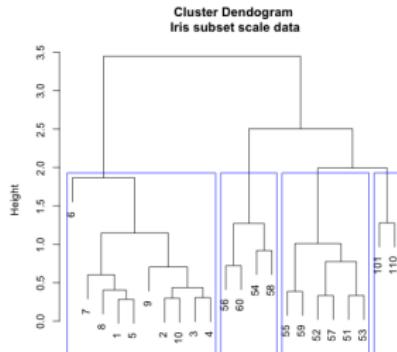


$D$   
hclust(\*, "ward.D")

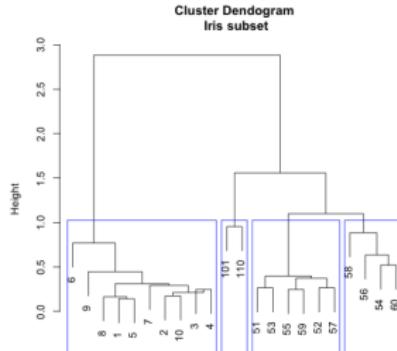
# Hierachical Clustering. Impact of Scaling. Illustration



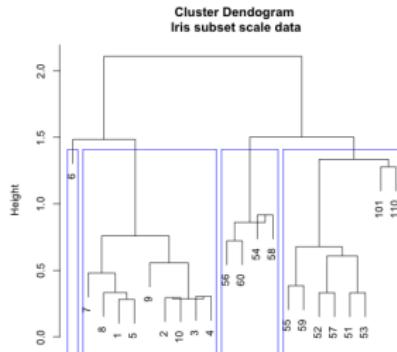
$D$   
`hclust(*, "average")`



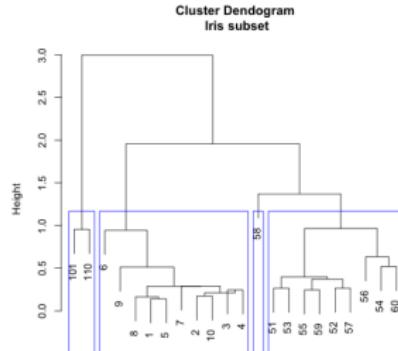
$D$   
`hclust(*, "average")`



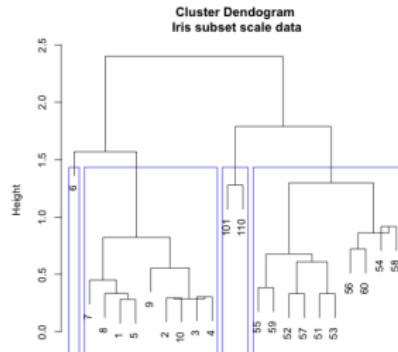
$D$   
`hclust(*, "centroid")`



$D$   
`hclust(*, "centroid")`

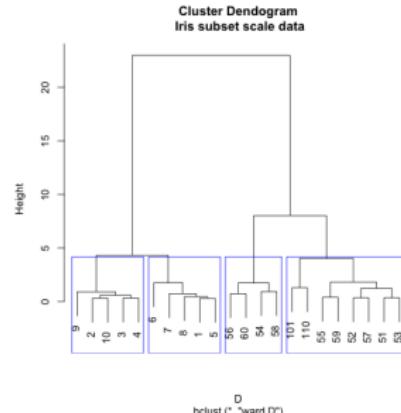
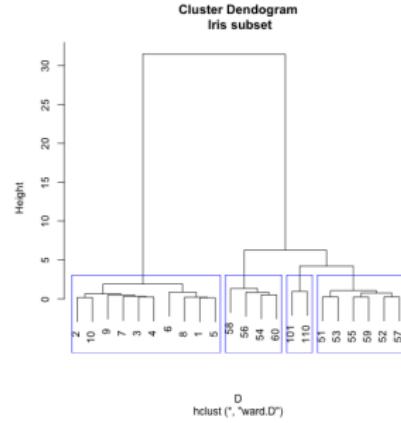
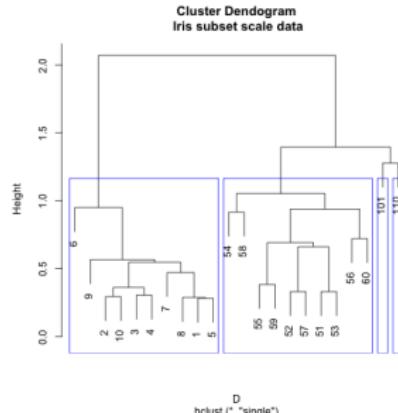
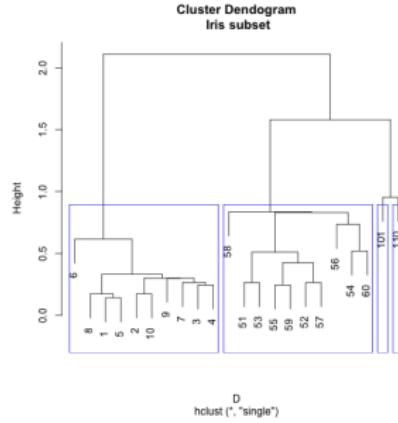
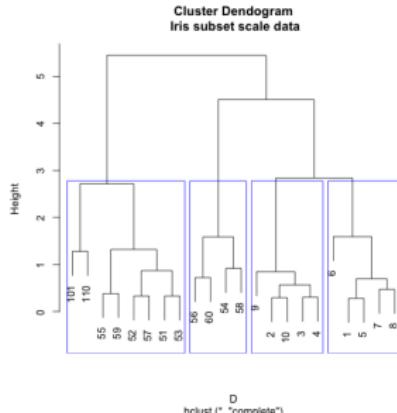
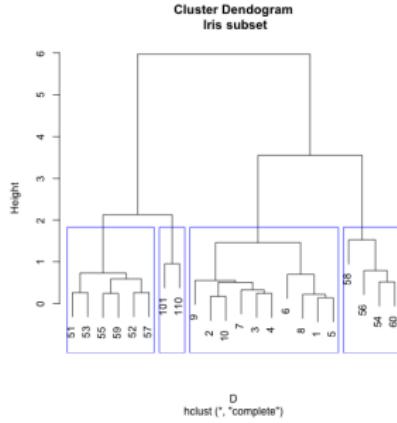


$D$   
`hclust(*, "median")`



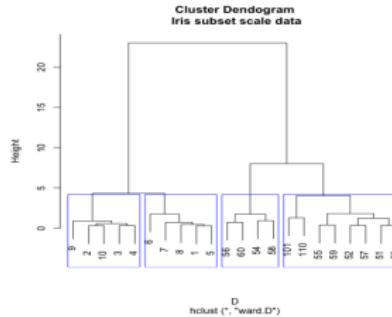
$D$   
`hclust(*, "median")`

# Hierachical Clustering. Impact of Scaling. Illustration



# Hierachical Clustering. R instructions.

```
# Tab : dataframe  
help(Hclust); D=dist(tab);  
HC=hclust(D,method="ward.D");  
plot(HC);  
rect.hclust(HC, k = 4);  
members <- cutree(HC, k = 4);
```



## Clustering method #2b -Top-down heuristic

**Parameters :** Threshold  $t$

**Initial step :** find the pair  $(x_i, x_j)$  having the maximal element in the distance matrix denoted  $d_M$

**Division step :** if  $d_M > t$  consider each of them as a center and affect remaining points to the closest center

**Iteration :** Iterate until  $d_M < t$  within each cluster.

# What about theory ?

- Pessimistic result (Kleinberg, 2003)

There is no clustering function satisfying scale invariance, richness and consistency.

- Optimistic result (von Luxburg, Ben-David, 2005)

Some stability for clustering can be guaranteed.