

Comment déterminer des profils cliniques homogènes

Journées OUTCOMEREA

4 décembre 2015

Sébastien BAILLY



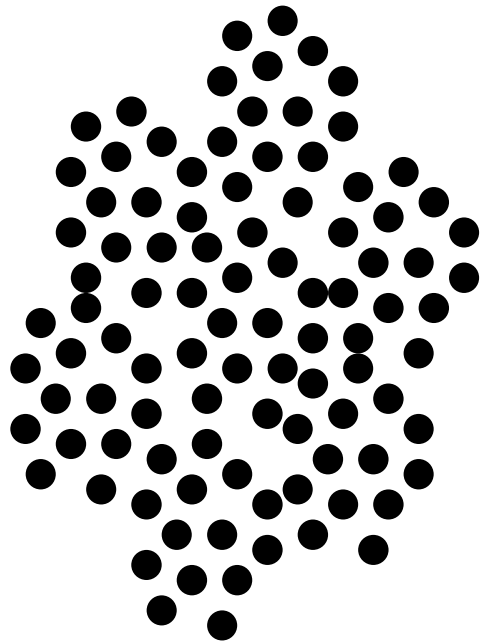
Introduction

Apprentissage non supervisé

Apprentissage supervisé

Données longitudinales

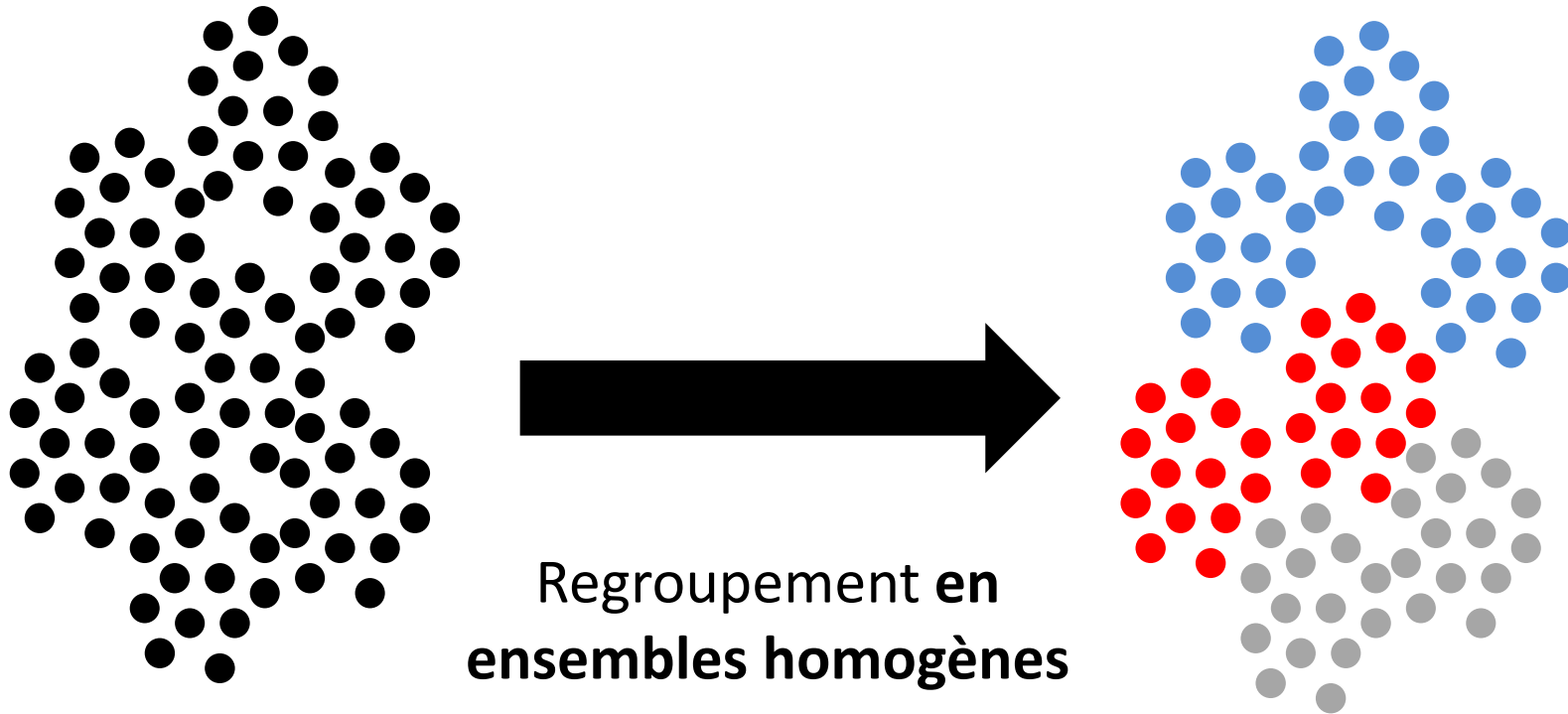
Conclusion



X patients

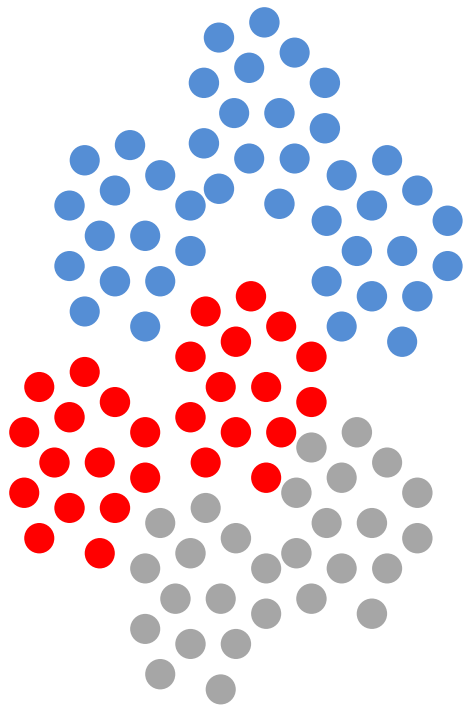


Variables cliniques
Biologiques
Pharmaceutiques
Démographiques
etc.



Classer les patients dans des classes
ou **clusters**

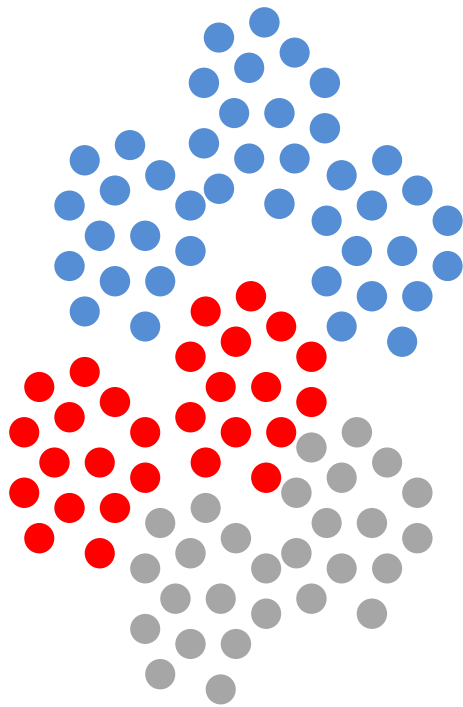
Cluster analysis



Minimiser la
distance entre deux
individus semblables



Maximiser la
distance entre deux
individus différents



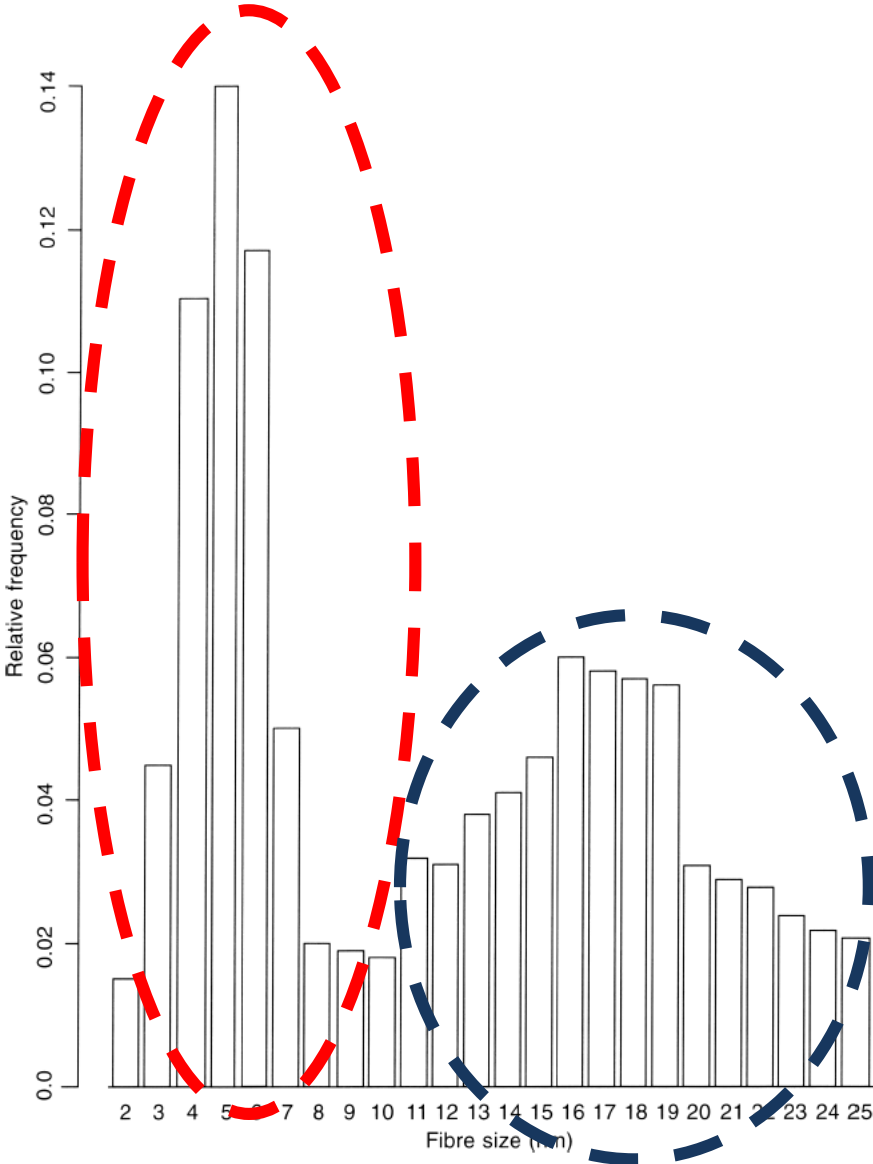
Minimiser la
distance entre deux
individus semblables



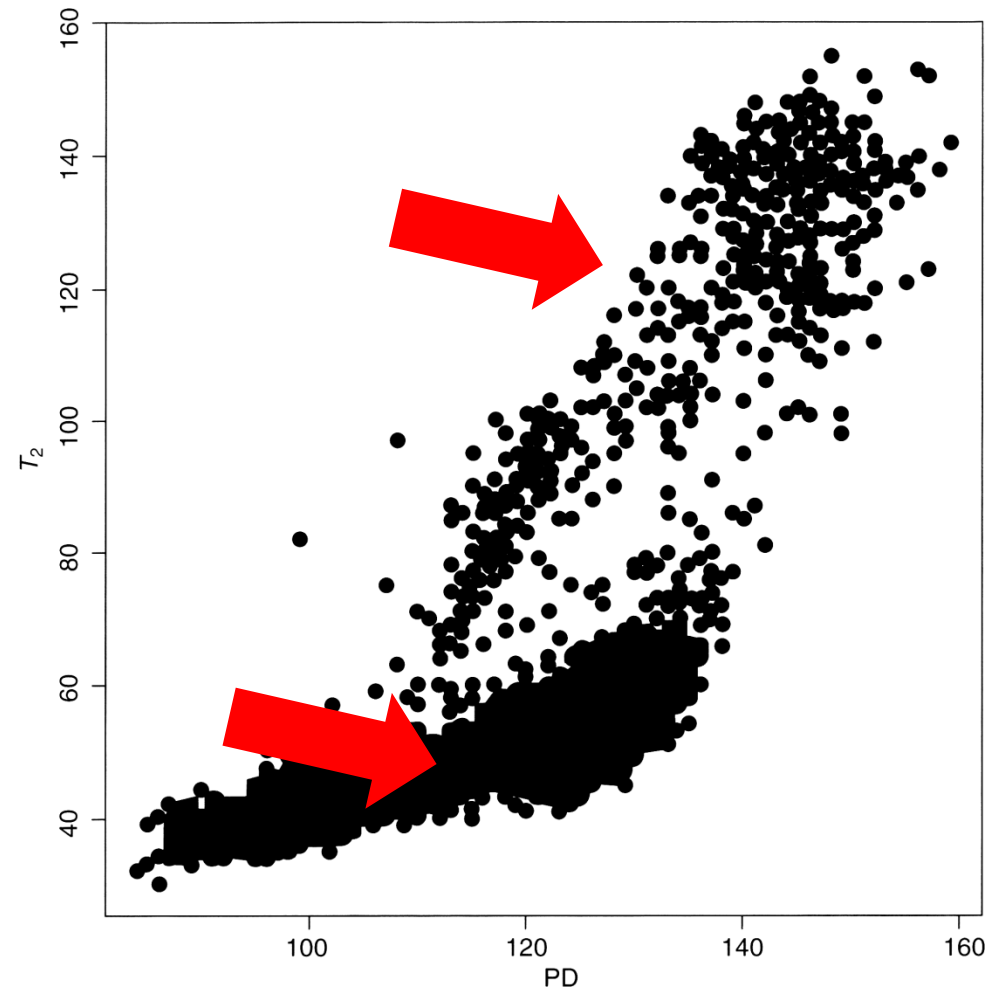
Maximiser la
distance entre deux
individus différents

Graphiques à deux dimensions

Histogrammes

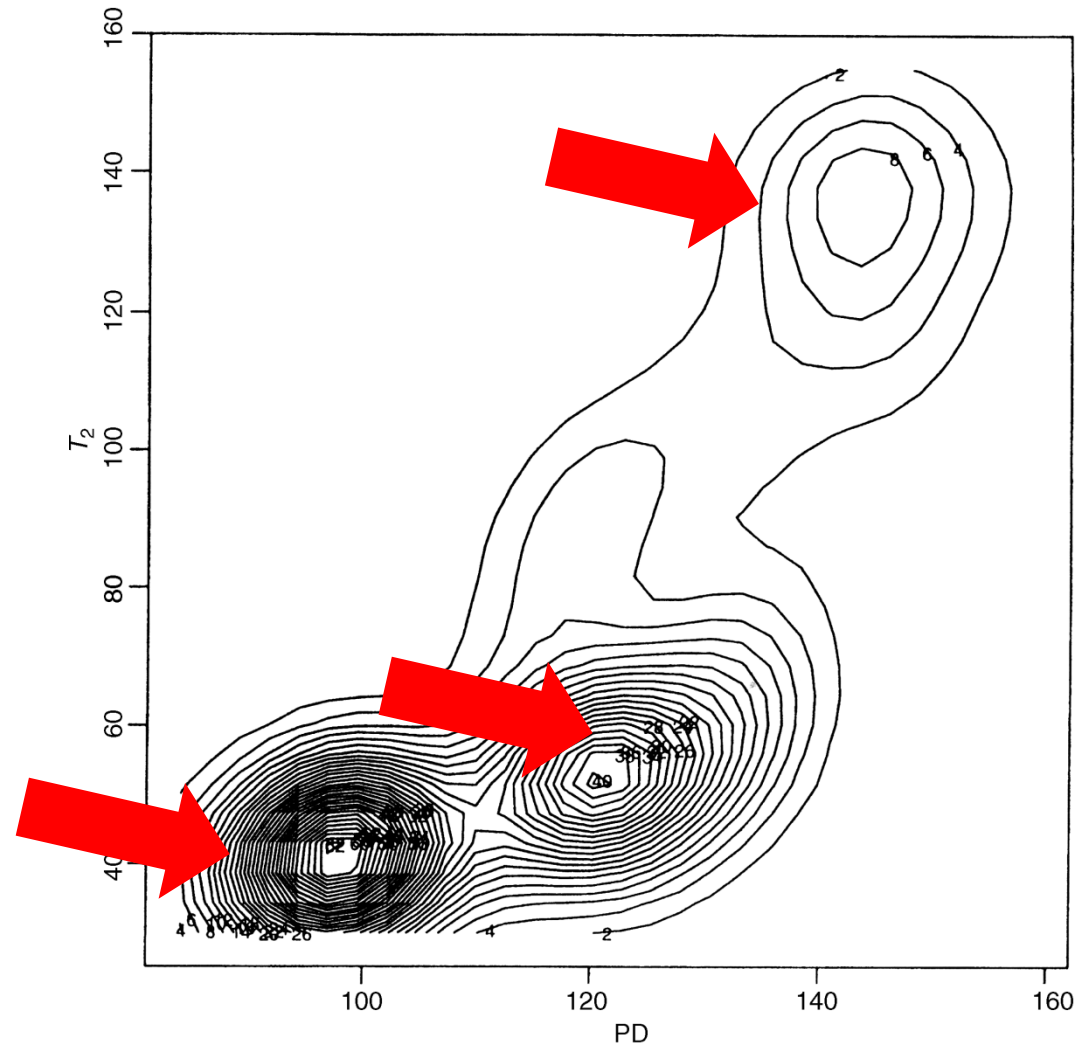


Graphiques à
deux dimensions
Nuages de points



Graphiques à deux dimensions

Courbes de densités



Augmentation du **nombre de variables** et du
nombre de patients



Approches plus complexes

Clusterisation
hiérarchique
ascendante

K-means

Analyses en
composantes
principales

Analyses de
classes
latentes

Analyses des
correspondances
multiples

Modèles de
mélanges

Réseaux de
neurones

Cartes de
Kohonen

Classification
naïve
bayésienne

Méthode
de Kernel

Classification
floue

Courbes de
densité

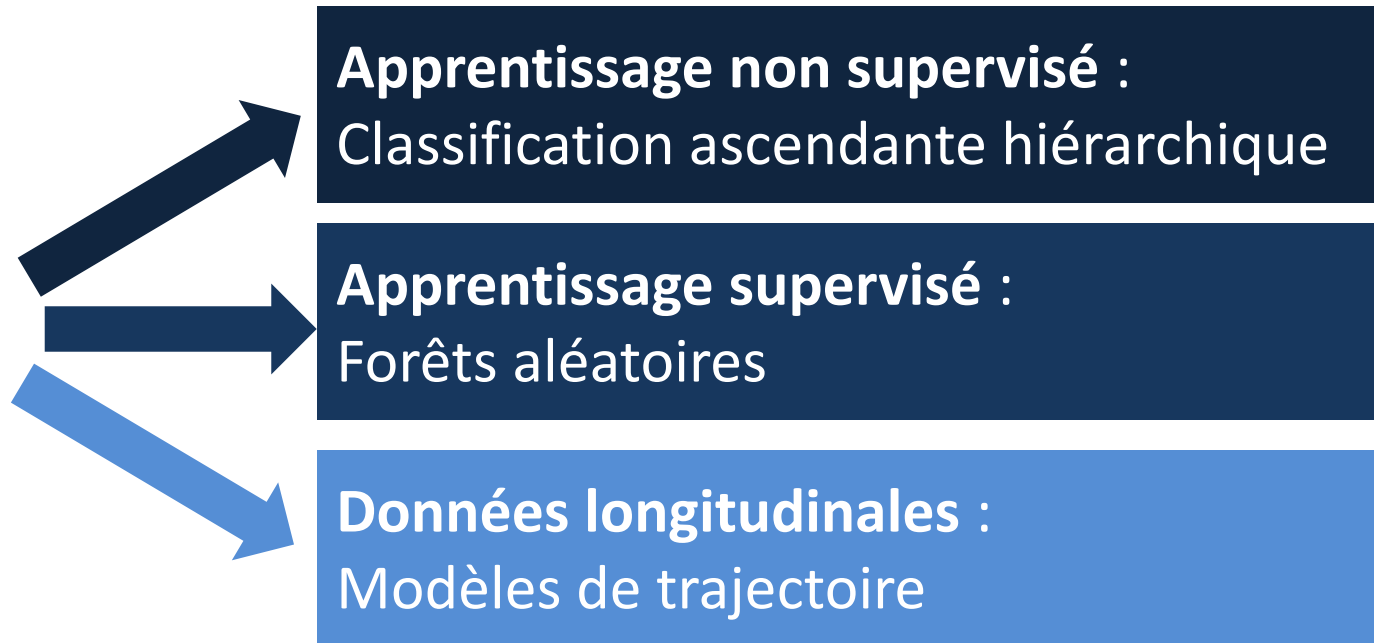
Modèles de
trajectoire

Augmentation du **nombre de variables** et du
nombre de **patients**



Approches plus complexes





Introduction



Apprentissage non supervisé

Apprentissage supervisé

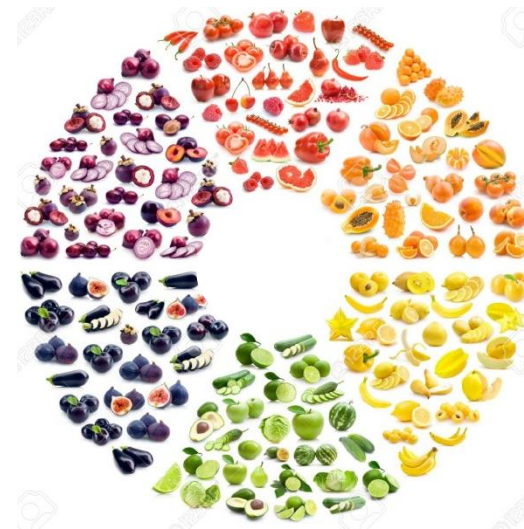
Données longitudinales

Conclusion

Recherche de classes
absence de données à expliquer



Recherche de classes absence de données à expliquer



Recherche de classes
absence de données à expliquer



Structuration des données après
regroupement



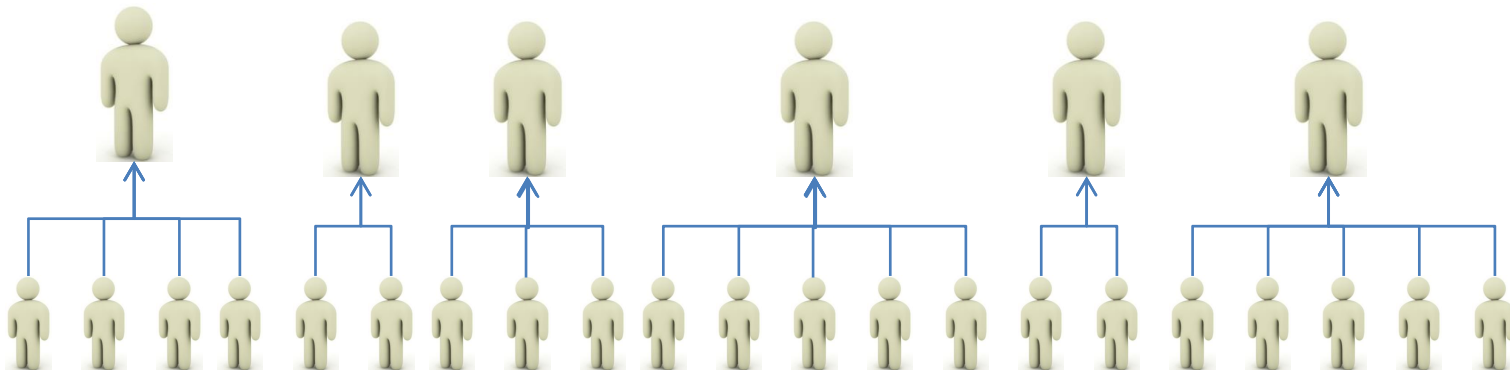
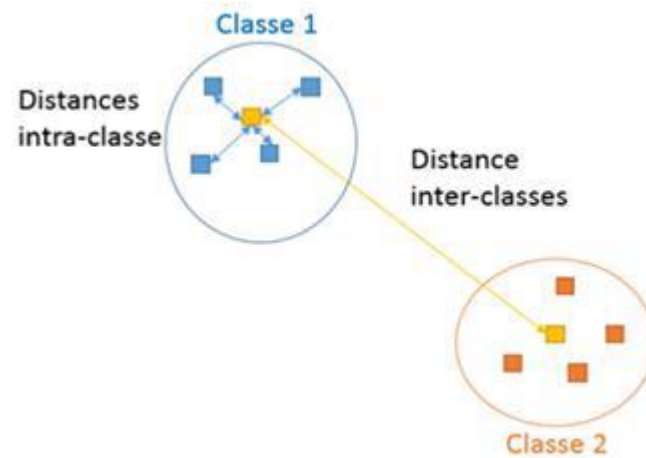
1^{ère} étape : chaque individu est son propre cluster

 Calcul de la matrice des distance des individus deux à deux



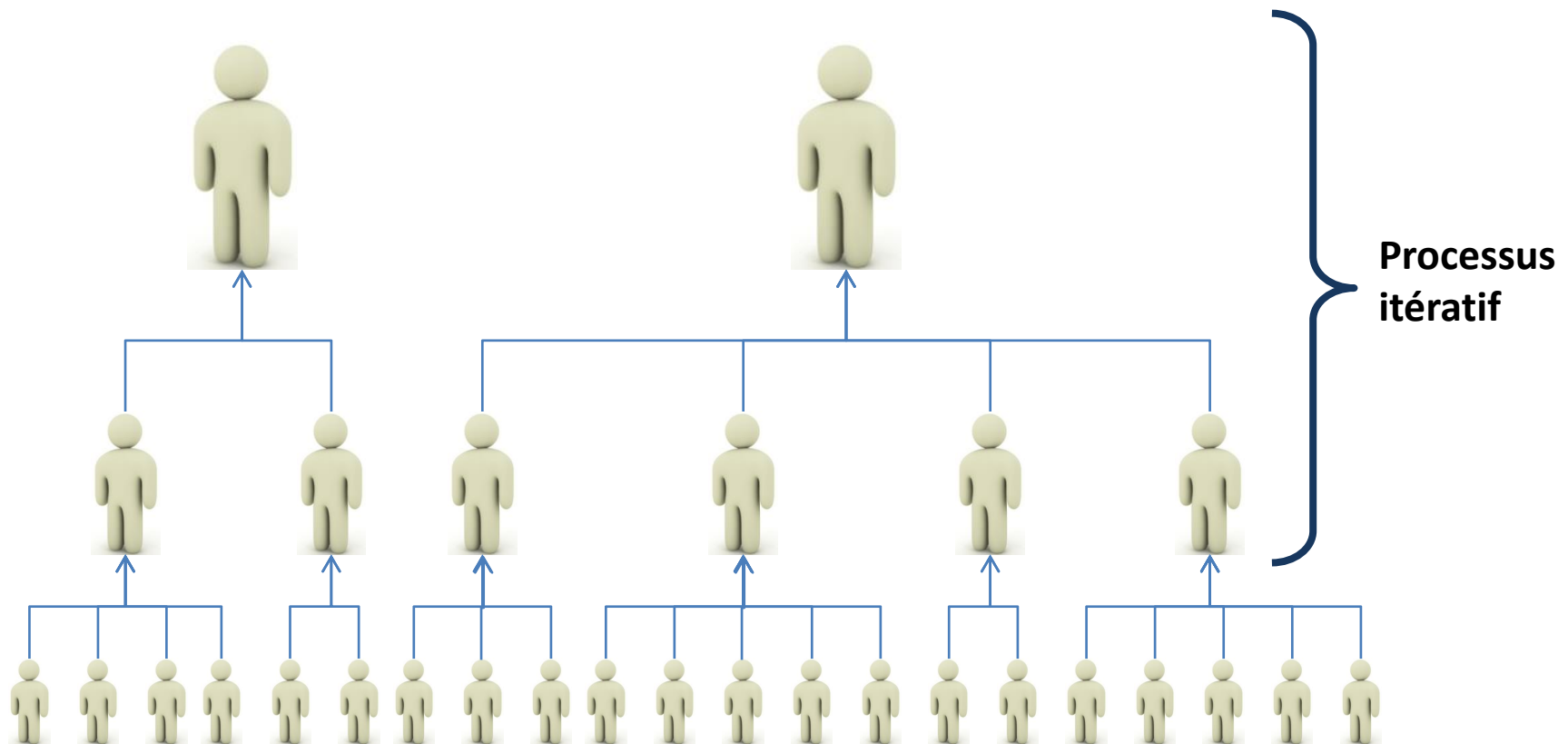
2^{ème} étape : Regroupement des éléments les plus proches

↳ Mise à jour de la matrice des distances en fonction des nouvelles classes

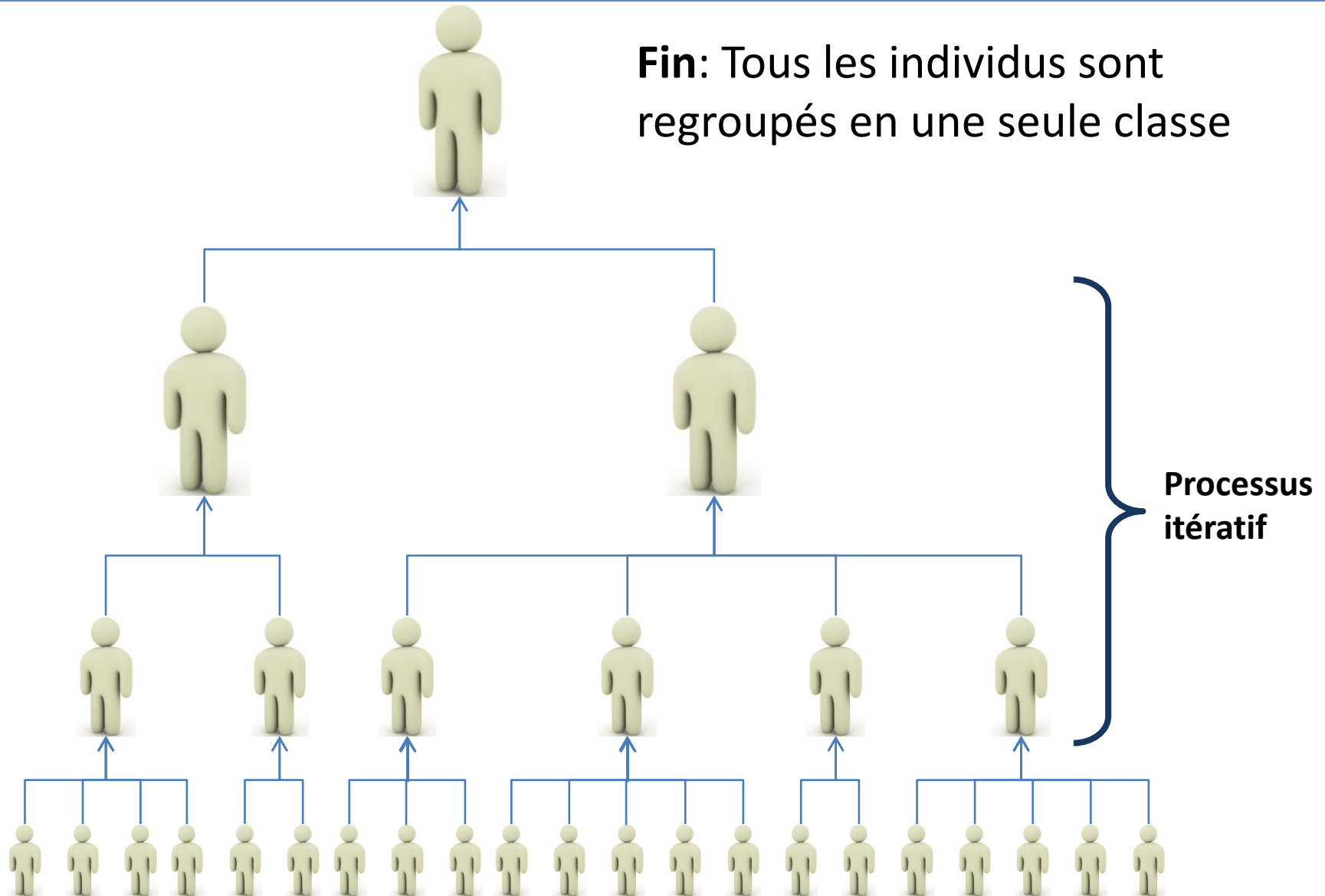


2^{ème} étape : Regroupement des éléments les plus proches

↳ Mise à jour de la matrice des distances en fonction des nouvelles classes



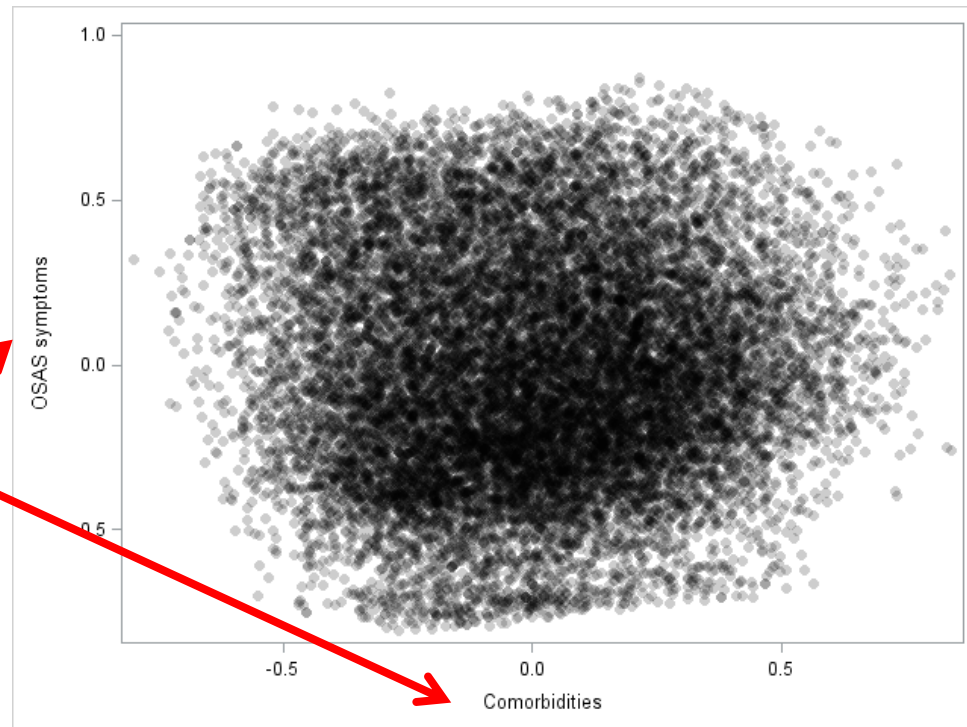
Fin: Tous les individus sont regroupés en une seule classe

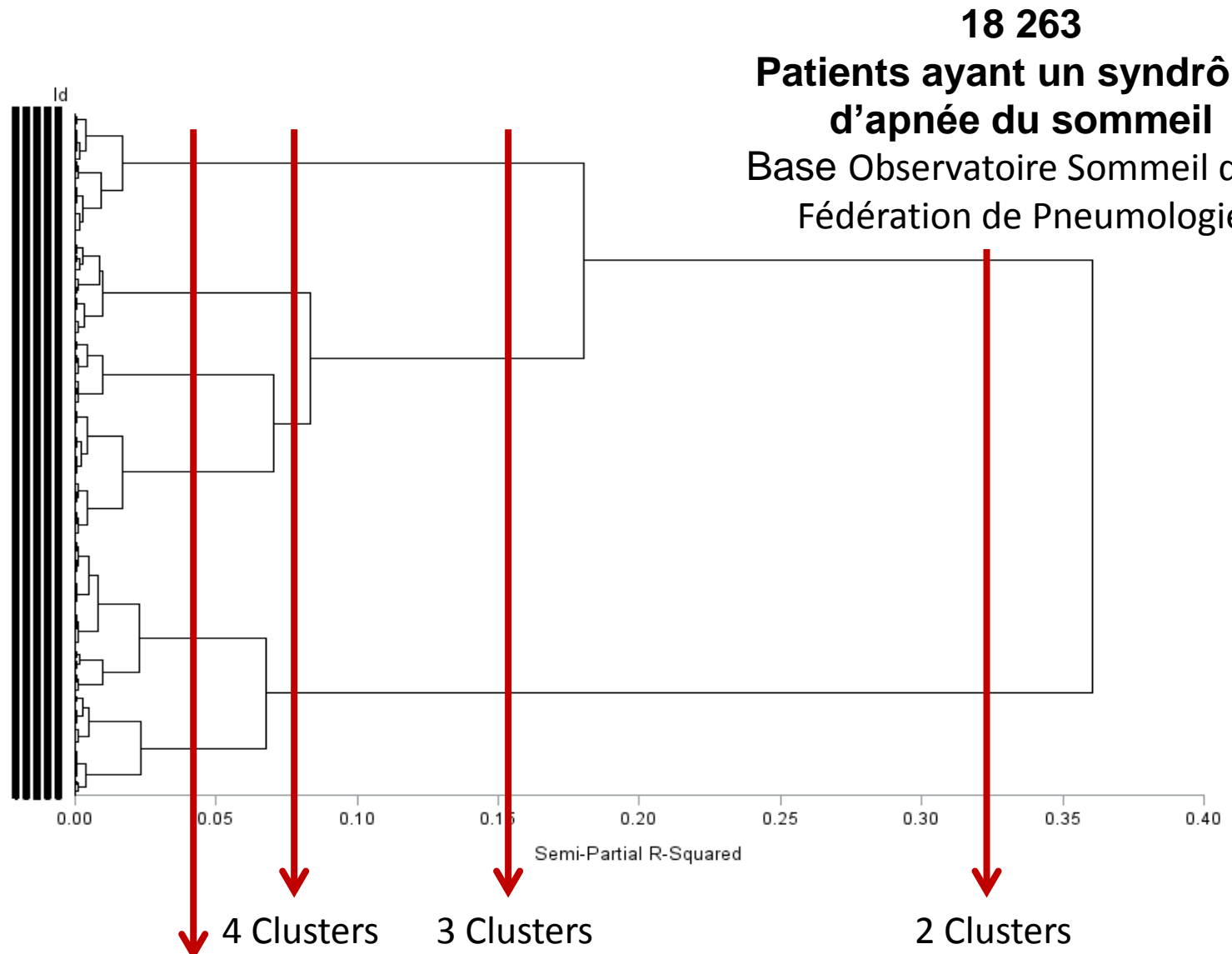


18 263

Patients ayant un syndrome d'apnée du sommeil
Base Observatoire Sommeil de la Fédération de Pneumologie

Réduction
de
l'information





Dendrogramme

Apprentissage non supervisé : classification ascendante hiérarchique

Cluster 1 (10%):

Jeunes (48a)

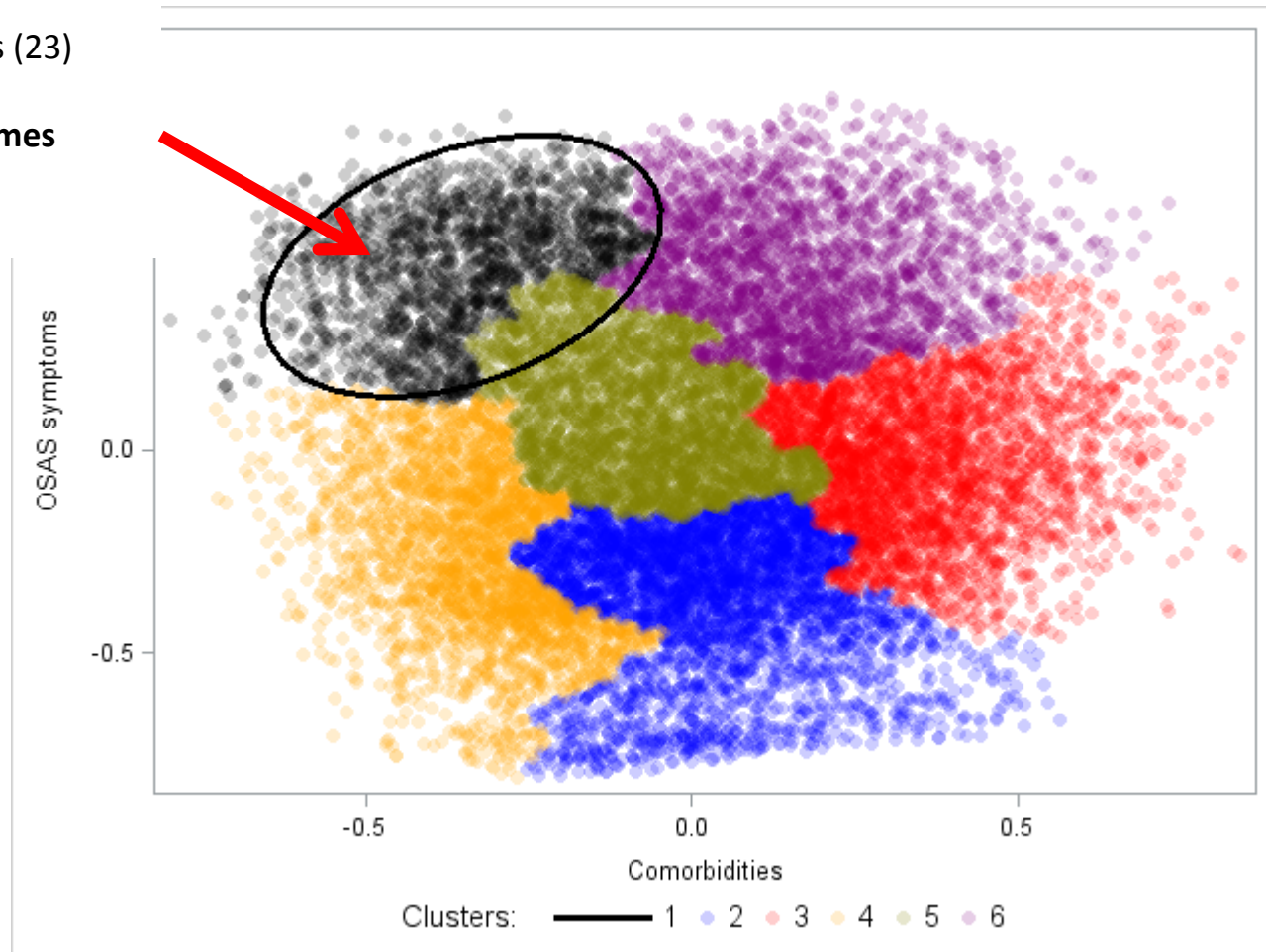
BMI bas (29)

AHI bas (32) – ODI bas (23)

Peu de comorbidités

Beaucoup de symptômes

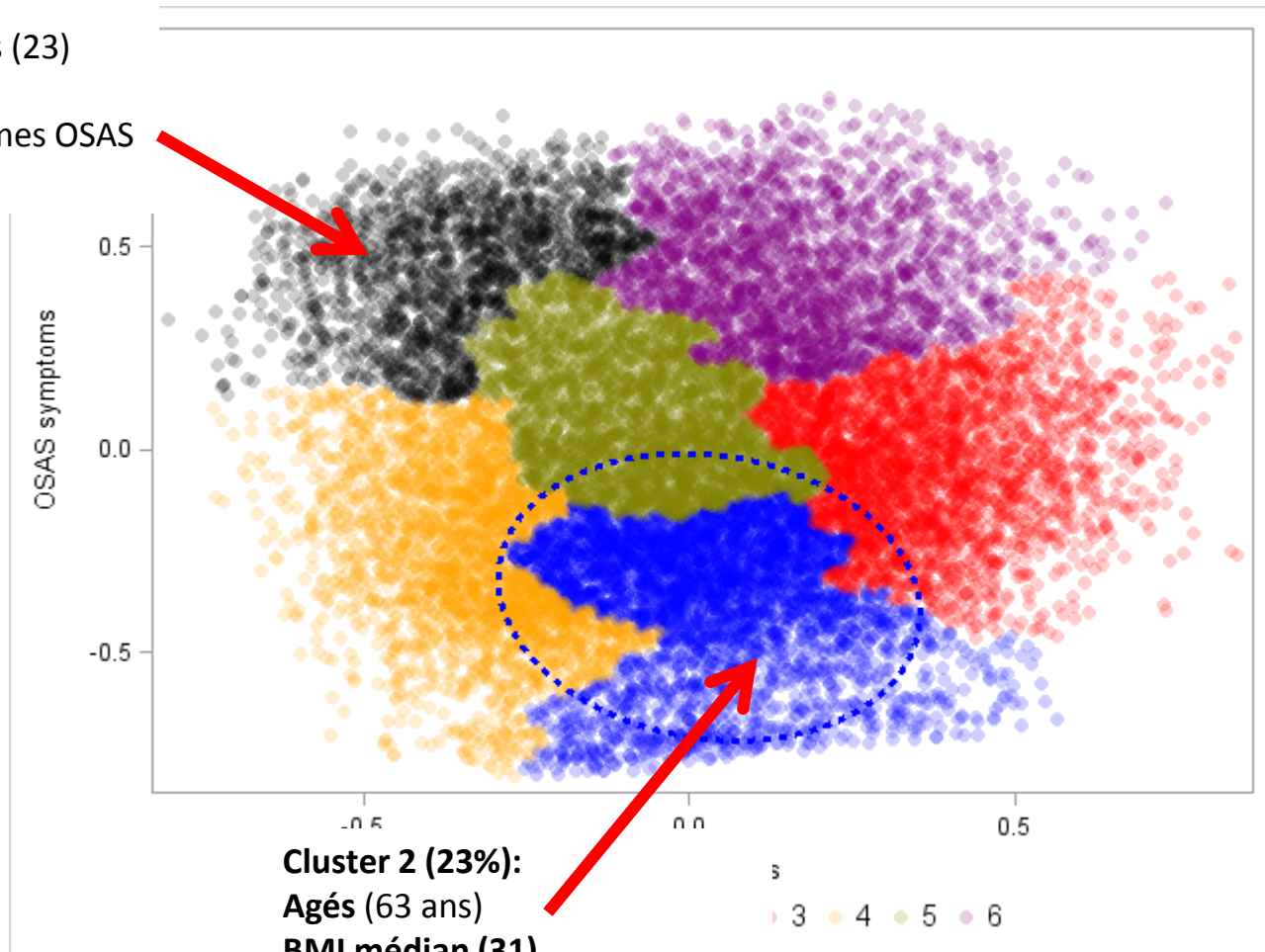
OSAS



Apprentissage non supervisé : classification ascendante hiérarchique

Cluster 1 (10%):

Jeunes (48a)
 BMI bas (29)
 AHI bas (32) – ODI bas (23)
 Peu de comorbidités
 Beaucoup de symptômes OSAS



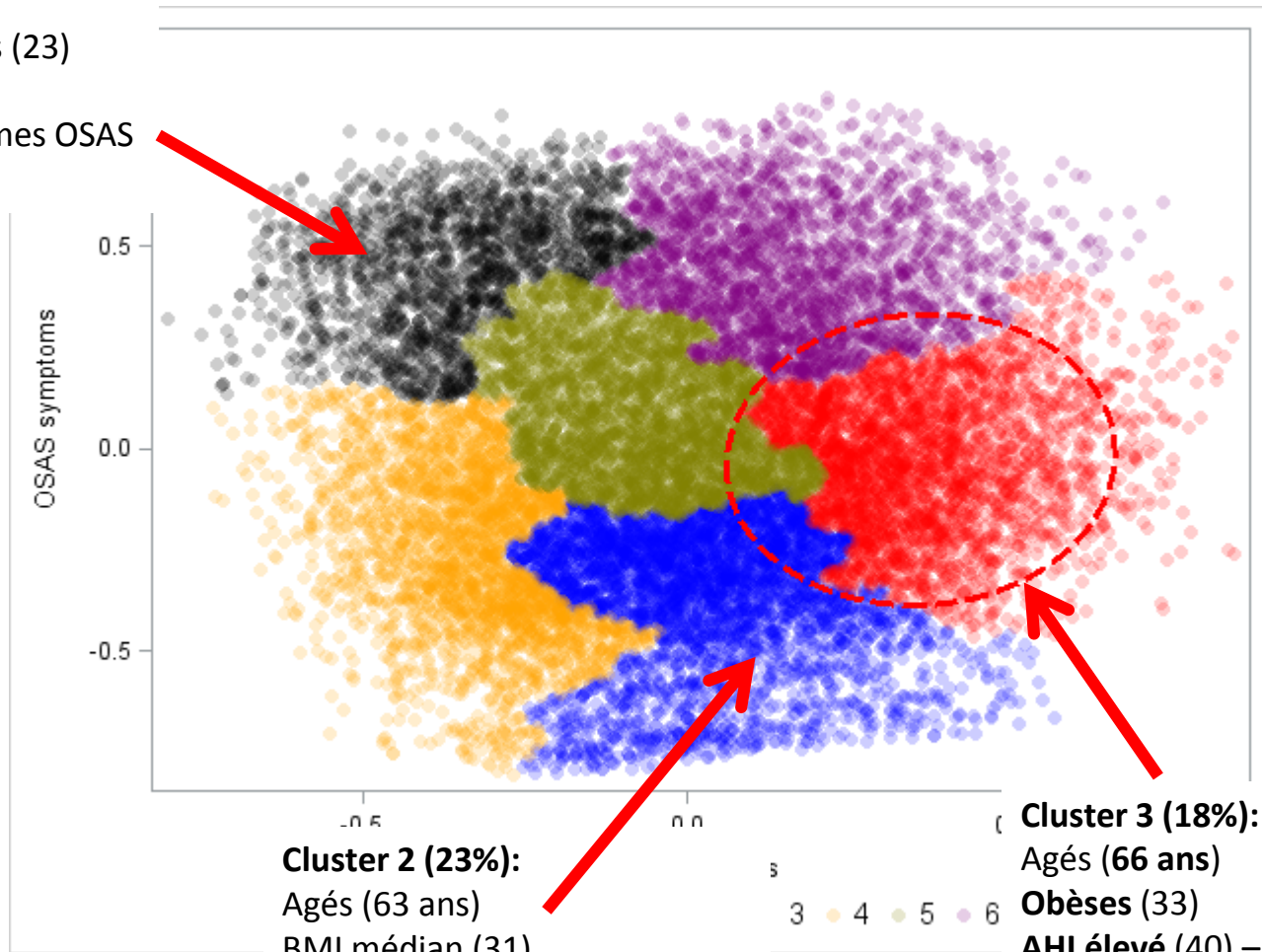
Cluster 2 (23%):

Agés (63 ans)
 BMI médian (31)
 AHI médian (34) – ODI bas (26)
 Peu de comorbidités
 Peu de symptômes OSAS

Apprentissage non supervisé : classification ascendante hiérarchique

Cluster 1 (10%):

Jeunes (48a)
 BMI bas (29)
 AHI bas (32) – ODI bas (23)
 Peu de comorbidités
 Beaucoup de symptômes OSAS



Cluster 2 (23%):

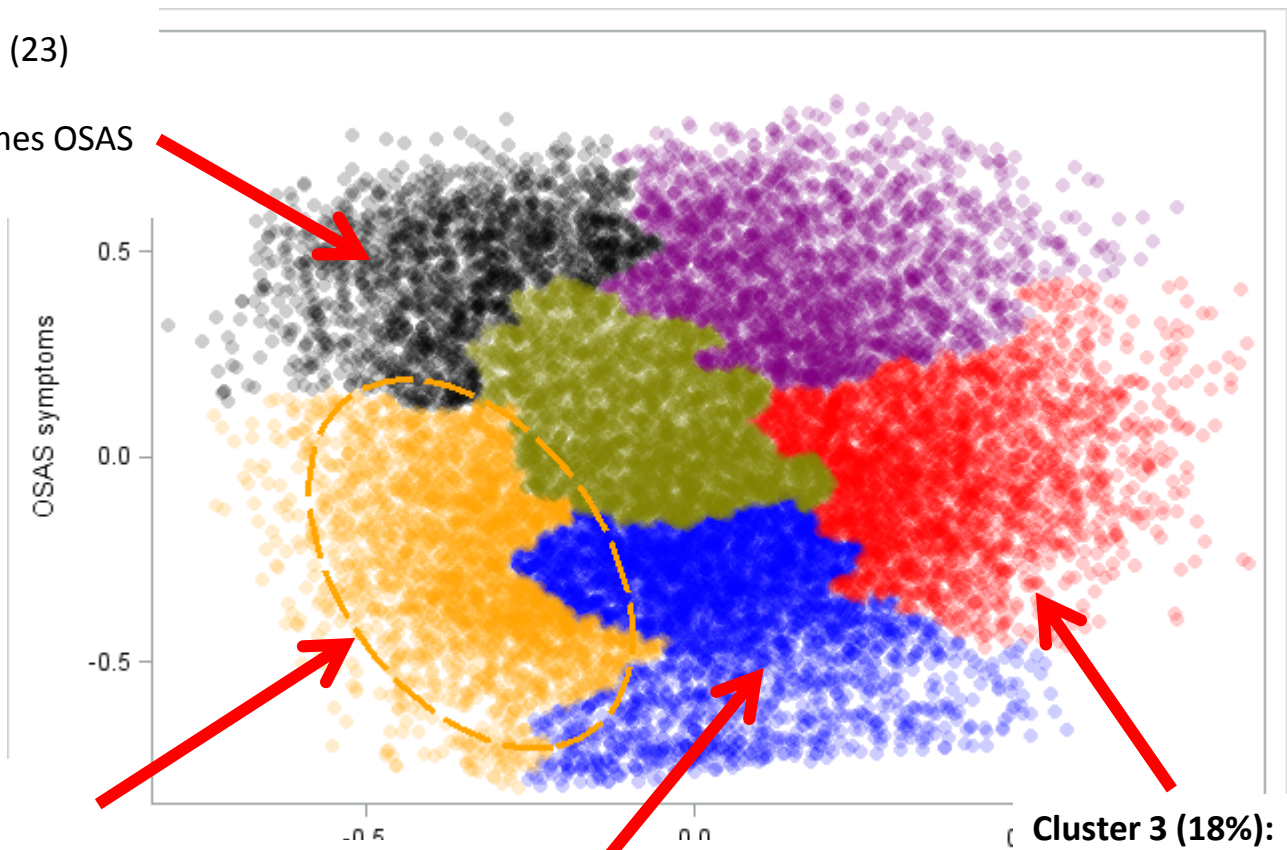
Agés (63 ans)
 BMI médian (31)
 AHI médian (34) – ODI bas (26)
 Peu de comorbidités
 Peu de symptômes OSAS

Cluster 3 (18%):

Agés (66 ans)
 Obèses (33)
 AHI élevé (40) – ODI élevé (33)
 Nombreuses comorbidités
 Peu de symptômes OSAS

Apprentissage non supervisé : classification ascendante hiérarchique

Cluster 1 (10%):
 Jeunes (48a)
 BMI bas (29)
 AHI bas (32) – ODI bas (23)
 Peu de comorbidités
 Beaucoup de symptômes OSAS



Cluster 4 (15%):
 Jeunes (49 ans)
 BMI bas (28)
 AHI bas (31) – ODI bas (21)
Peu de comorbidités
Peu de symptômes OSAS

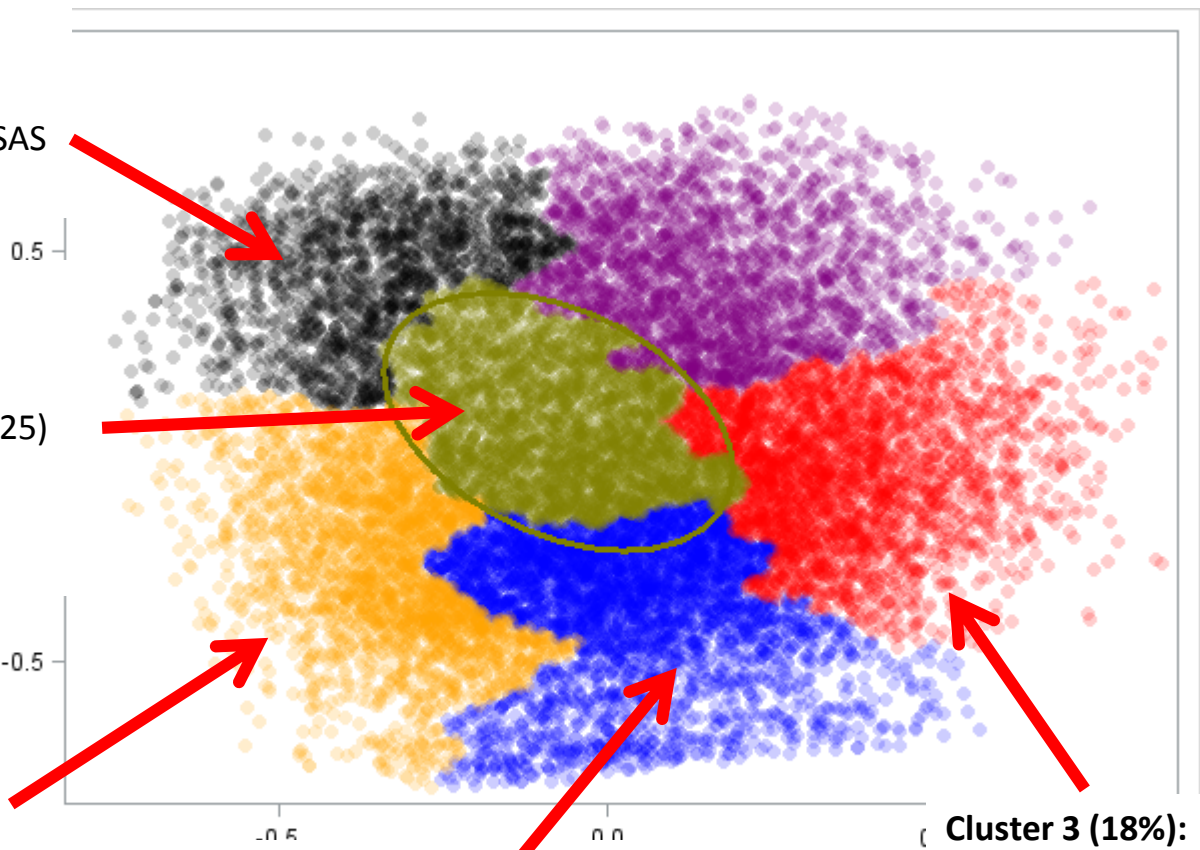
Cluster 2 (23%):
 Agés (63 ans)
 BMI médian (31)
 AHI médian (34) – ODI bas (26)
 Peu de comorbidités
 Peu de symptômes OSAS



Cluster 3 (18%):
 Agés (66 ans)
 Obèses (33)
 AHI élevé (40) – ODI élevé (33)
 Nombreuses comorbidités
 Peu de symptômes OSAS

Cluster 1 (10%):

Jeunes (48a)
BMI bas (29)
AHI bas (32) – ODI bas (23)
Peu de comorbidités
Beaucoup de symptômes OSAS



Cluster 5 (19%):

Age moyen (56 ans)
BMI moyen (31)
AHI moyen (34) – ODI moyen (25)
Peu de comorbidités
Peu de symptômes OSAS

Cluster 4 (15%):

Jeunes (49 ans)
BMI bas (28)
AHI bas (31) – ODI bas (21)
Peu de comorbidités
Peu de symptômes OSAS

Cluster 2 (23%):

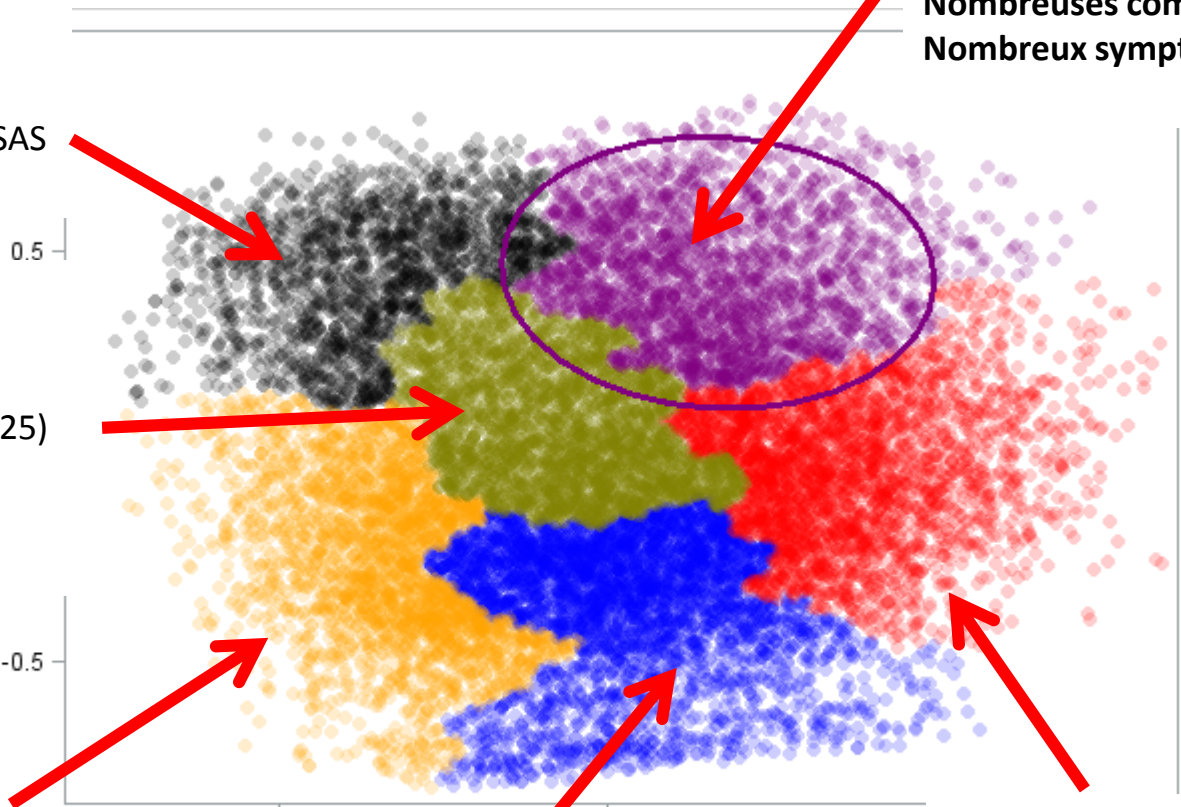
Agés (63 ans)
BMI médian (31)
AHI médian (34) – ODI bas (26)
Peu de comorbidités
Peu de symptômes OSAS

Cluster 3 (18%):

Agés (66 ans)
Obèses (33)
AHI élevé (40) – ODI élevé (33)
Nombreuses comorbidités
Peu de symptômes OSAS

Apprentissage non supervisé : classification ascendante hiérarchique

Cluster 1 (10%):
Jeunes (48a)
BMI bas (29)
AHI bas (32) – ODI bas (23)
Peu de comorbidités
Beaucoup de symptômes OSAS



Cluster 6 (14%):
Age moyen (60 ans)
Obèses (33)
AHI élevé (39) – ODI élevé (31)
Nombreuses comorbidités
Nombreux symptômes OSAS

Cluster 5 (19%):
Age moyen (56 ans)
BMI moyen (31)
AHI moyen (34) – ODI moyen (25)
Peu de comorbidités
Peu de symptômes OSAS

Cluster 4 (15%):
Jeunes (49 ans)
BMI bas (28)
AHI bas (31) – ODI bas (21)
Peu de comorbidités
Peu de symptômes OSAS

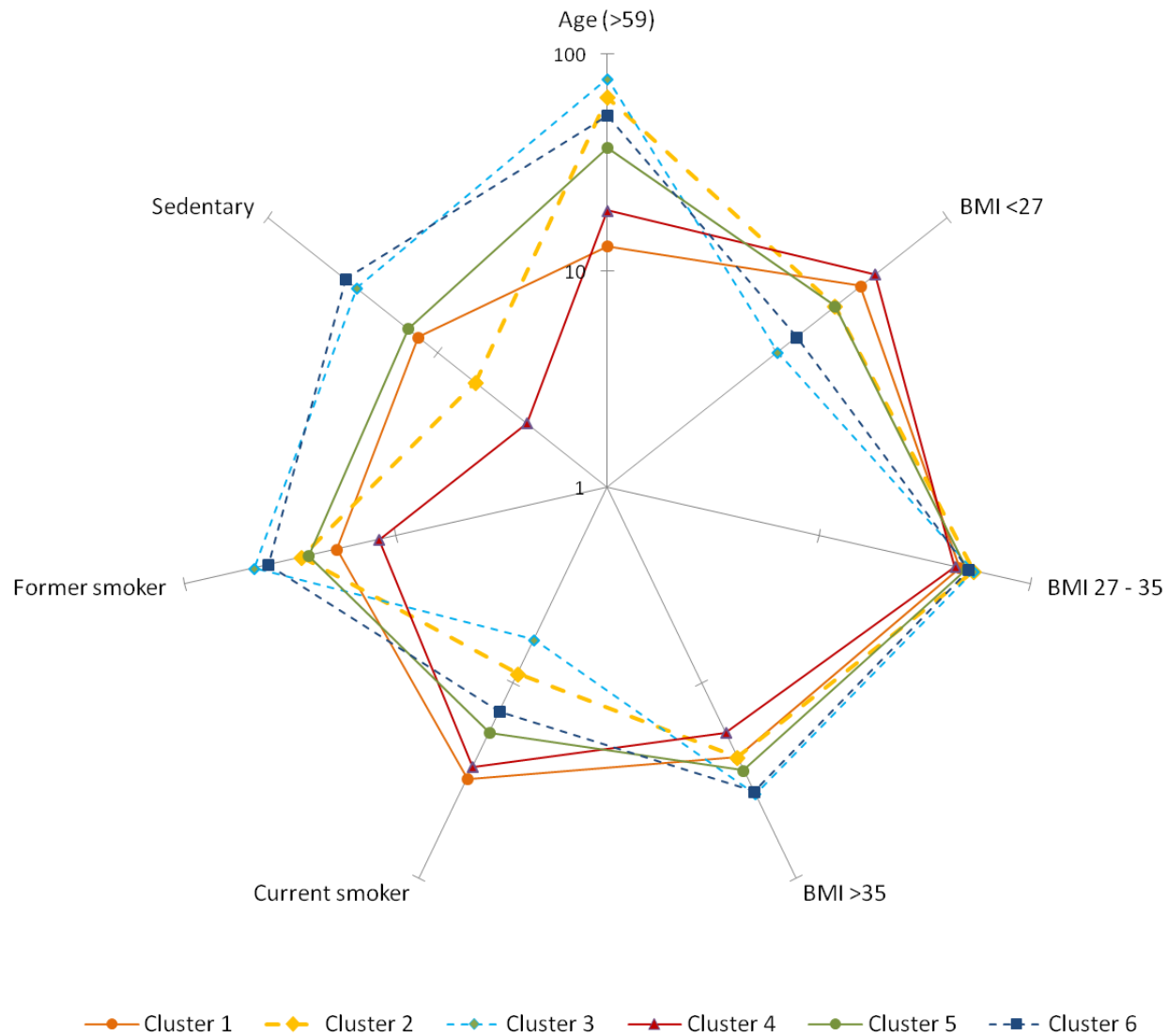
Cluster 2 (23%):
Agés (63 ans)
BMI médian (31)
AHI médian (34) – ODI bas (26)
Peu de comorbidités
Peu de symptômes OSAS

Cluster 3 (18%):
Agés (66 ans)
Obèses (33)
AHI élevé (40) – ODI élevé (33)
Nombreuses comorbidités
Peu de symptômes OSAS

Apprentissage non supervisé : classification ascendante hiérarchique

	Age (years)	BMI (kg/m ²)	AHI (/h)	ODI (/h)	Epworth scale	Co-morbidities	OSAS symptoms
Cluster 1	Young (48)	Low (29)	Low (31.6)	Low (23)	High (12)	Few or no	Many
Cluster 2	Oldest (63)	Median (31)	Median (34)	Median (26)	Low (8)	Few or no	Few or no
Cluster 3	Oldest (66)	Obese (33)	High (40)	High (33)	Low (9)	Many	Few or no
Cluster 4	Young (49)	Low (28)	Low (31)	Low (21)	Median (10)	Few or no	Few or no
Cluster 5	Middle age (56)	Median (31)	Median (34)	Median (25)	High (11)	Few or no	Few or no
Cluster 6	Middle age (60)	Obese (33)	High (39)	High (31)	High (11)	Many	Many

➔ Facteurs environnementaux



- ➔ Méthode basée sur les distance entre les individus
- ➔ Qualification et interprétation simple des classes créées
- ➔ Méthode descriptive
- ➔ Dépend du jeu de données : nécessité de faire une validation externe

Introduction

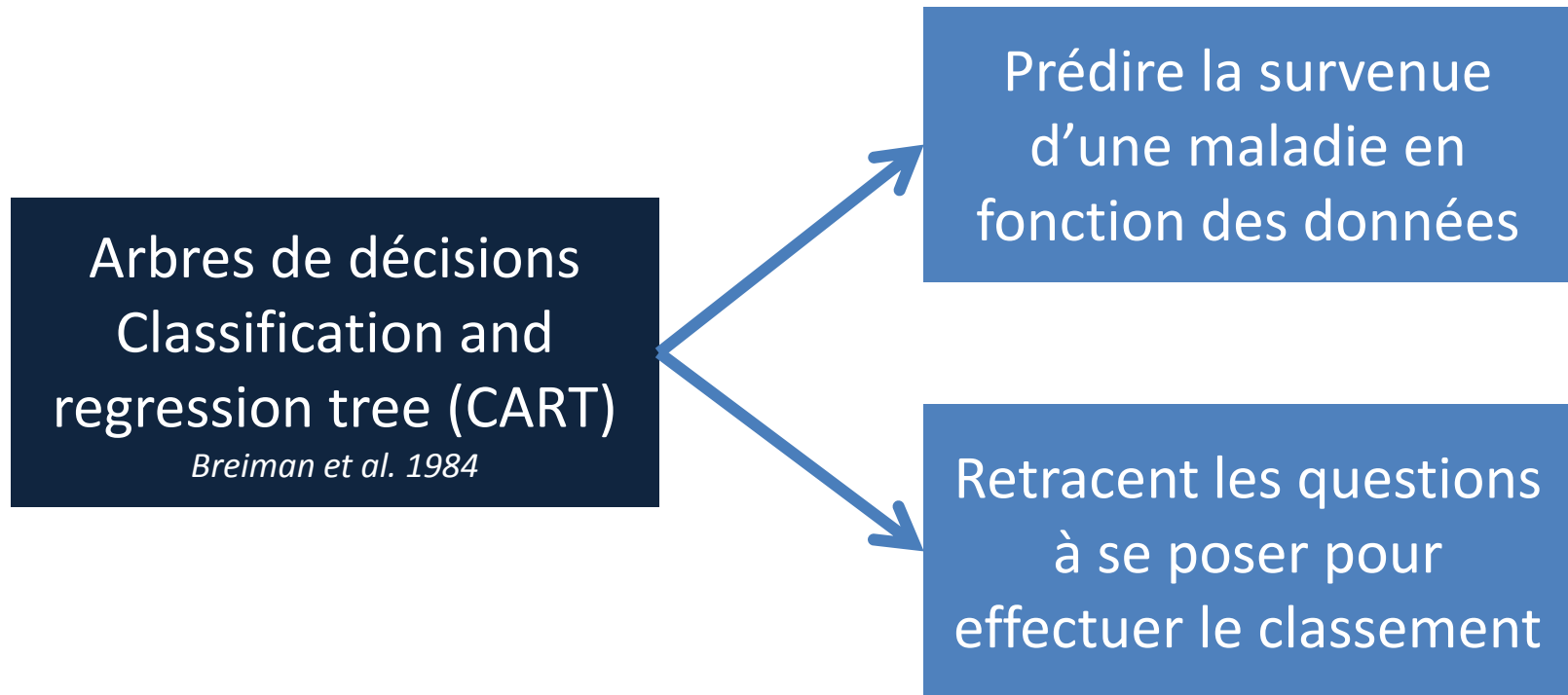
Apprentissage non supervisé

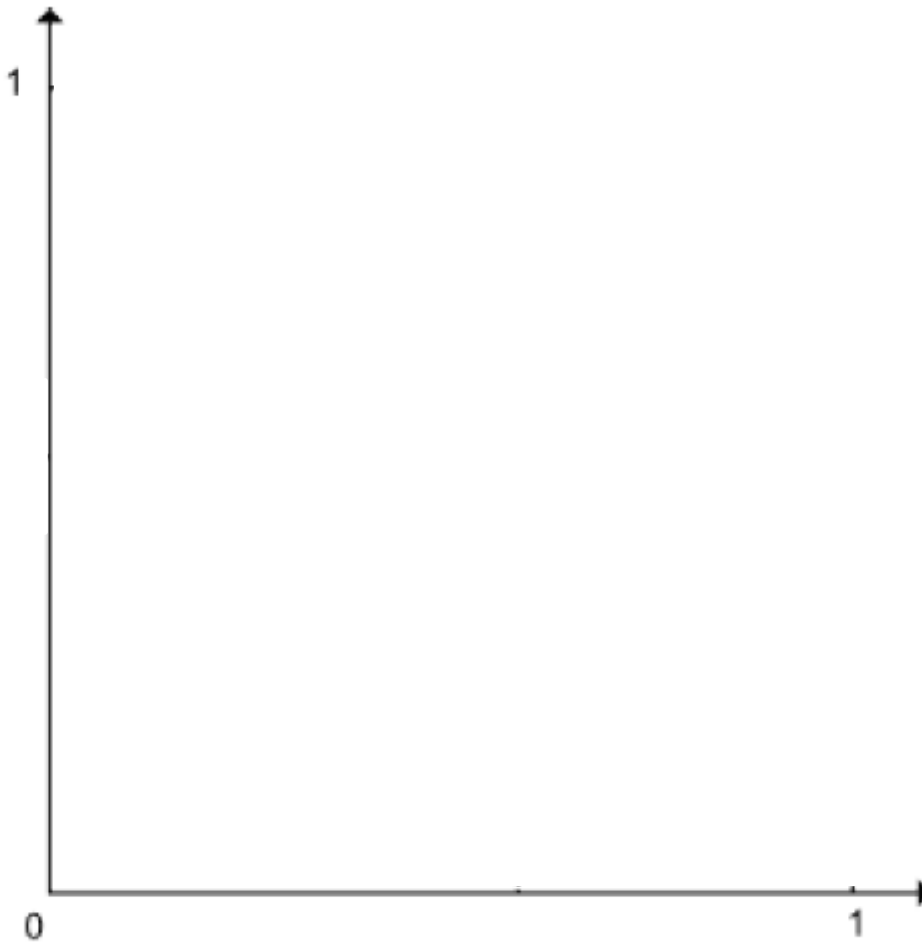


Apprentissage supervisé

Données longitudinales

Conclusion





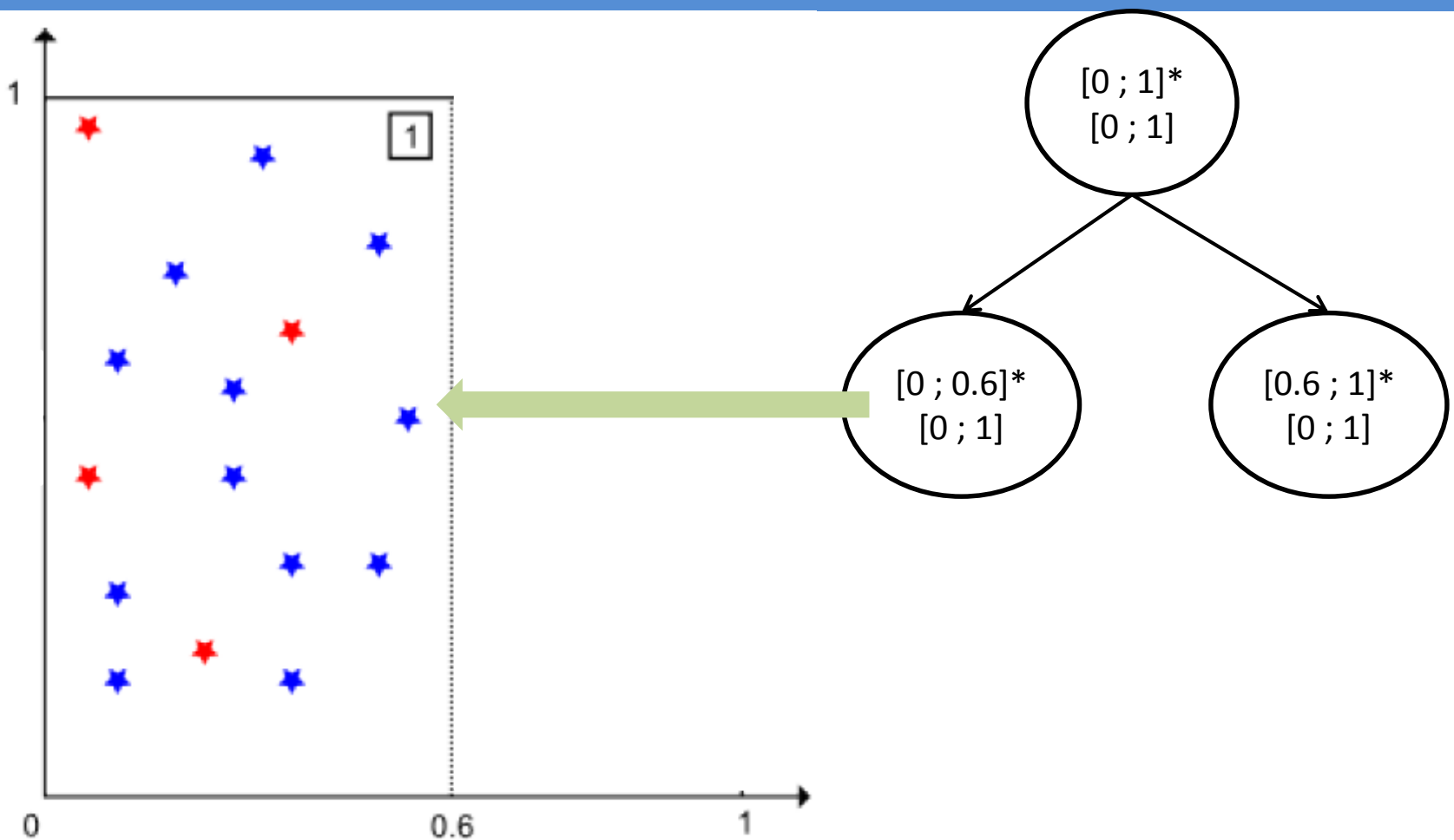
$[0 ; 1]^*$
 $[0 ; 1]$

On cherche à prédire une variable binaire



En fonction de deux variables
quantitatives
comprises entre 0 et 1

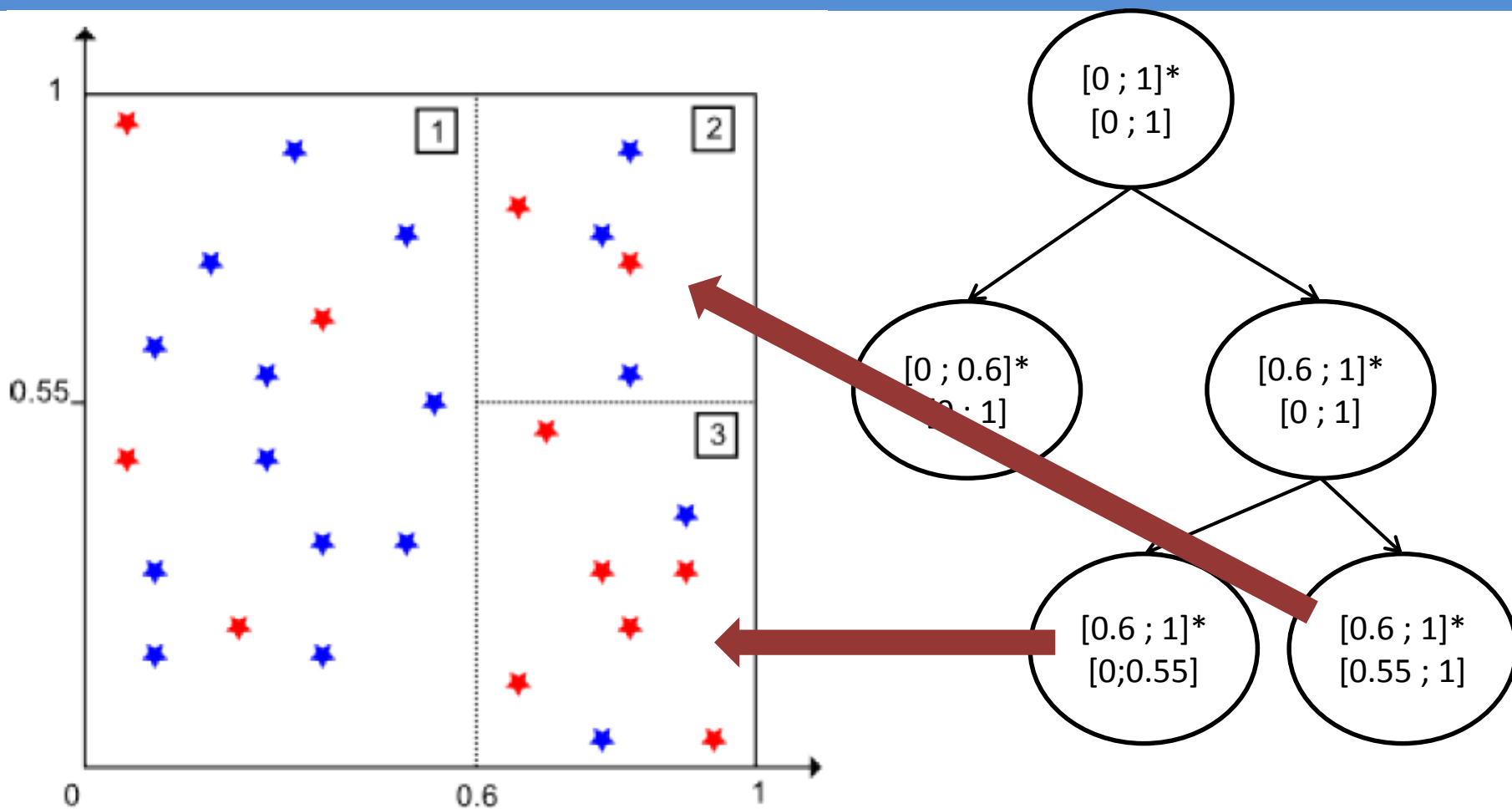
Apprentissage supervisé : Forêts aléatoires – Random forest



On cherche à prédire une variable binaire {
 \star (blue)
 \star (red)

En fonction de deux variables quantitatives comprises entre 0 et 1

Apprentissage supervisé : Forêts aléatoires – Random forest

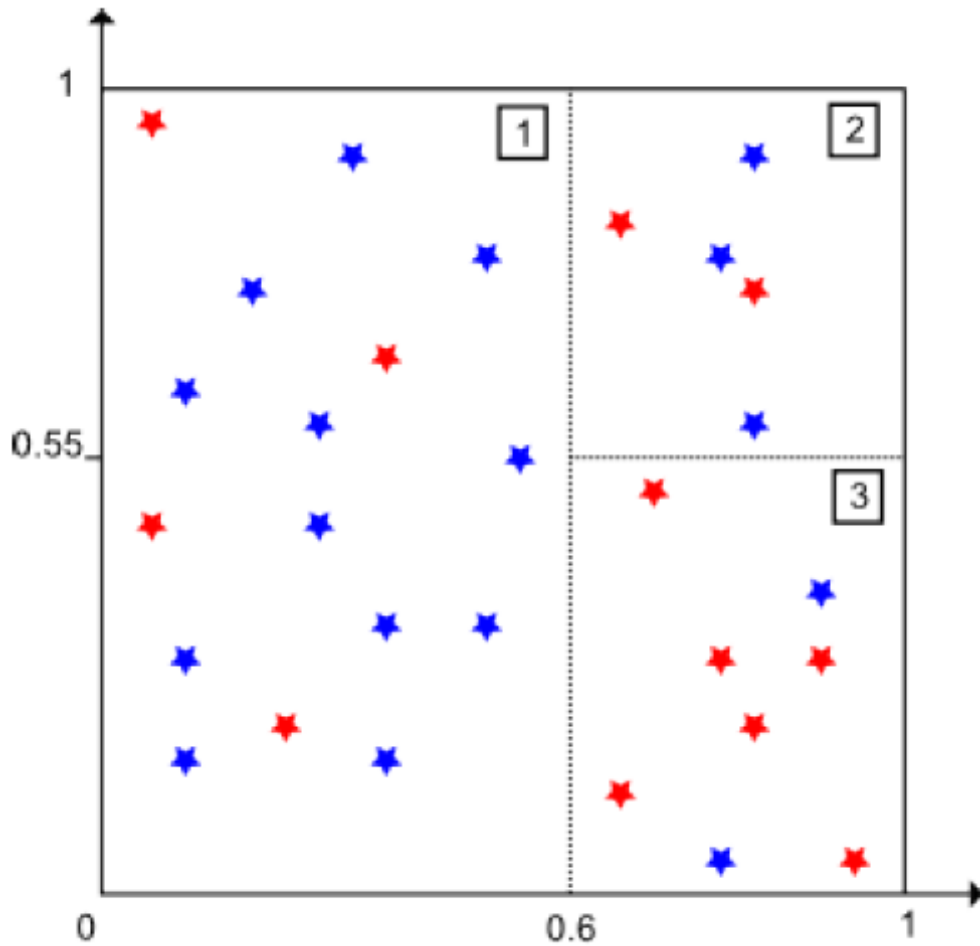


On cherche à prédire une variable binaire

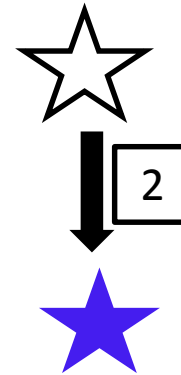
$\left\{ \begin{array}{l} \star \\ \star \end{array} \right.$

En fonction de deux variables
quantitatives
comprises entre 0 et 1

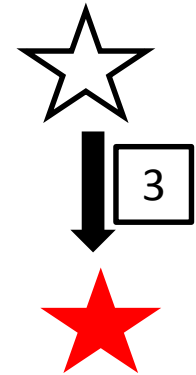
Apprentissage supervisé : Forêts aléatoires – Random forest









$0.72 * 0.75$



$0.70 * 0.20$



-  L'arbre peut être directement traduit en règles claires et interprétables
-  Sélection automatique des variables pertinentes
Robuste par rapport aux variables redondantes
-  Robuste par rapport aux variables aberrantes
Possibilité de prendre en compte les données manquantes
-  Traite rapidement de très grandes bases de données
-  Possibilité d'intervenir dans la construction de l'arbre
-  Sont fortement dépendant de l'échantillon initial

Random forest : Breiman 2001

Introduction d'une partie aléatoire dans la construction des arbres de décision

Classification and regression Tree



Un seul arbre vote pour classer un nouveau patient malade ou non malade



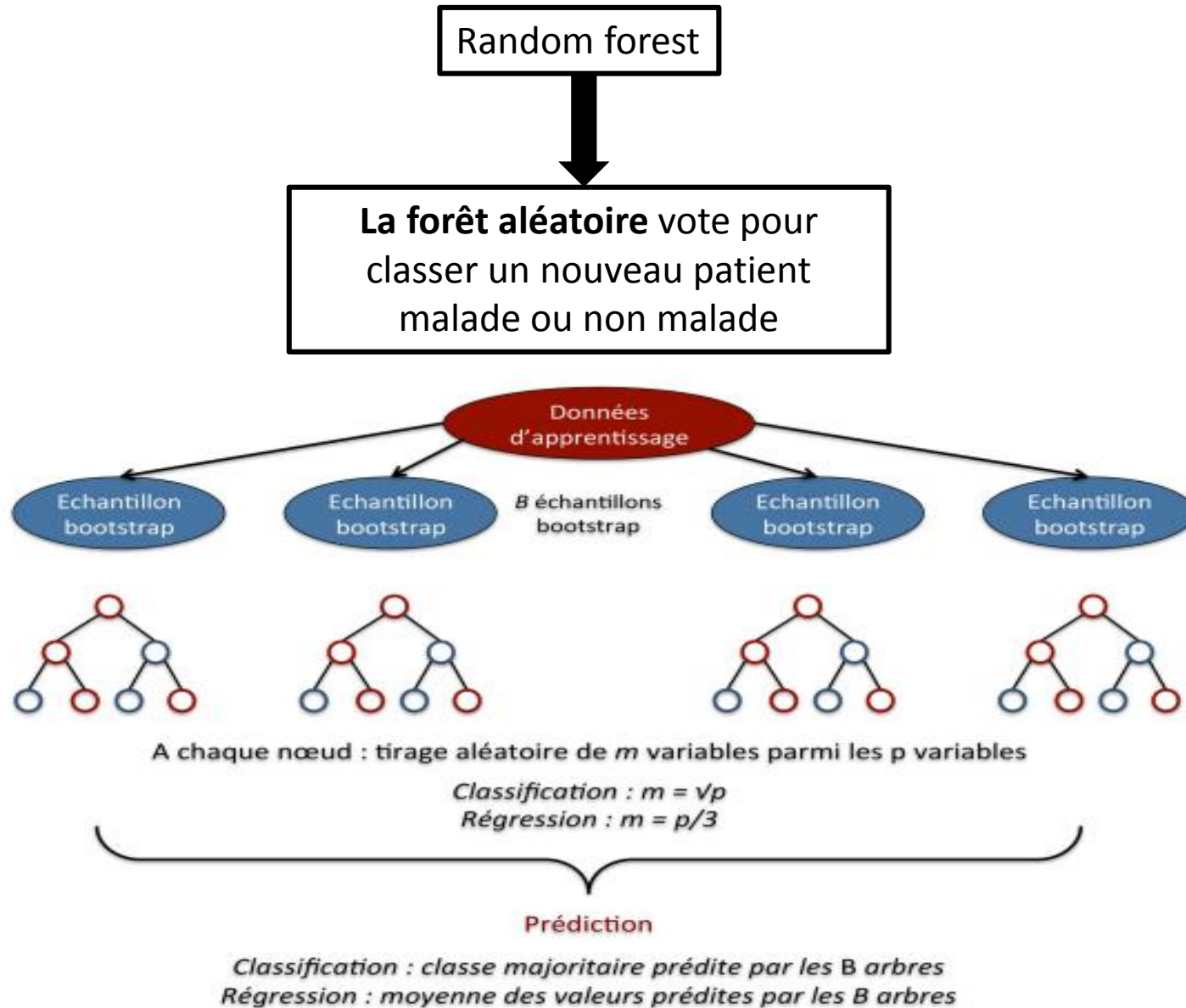
Random forest



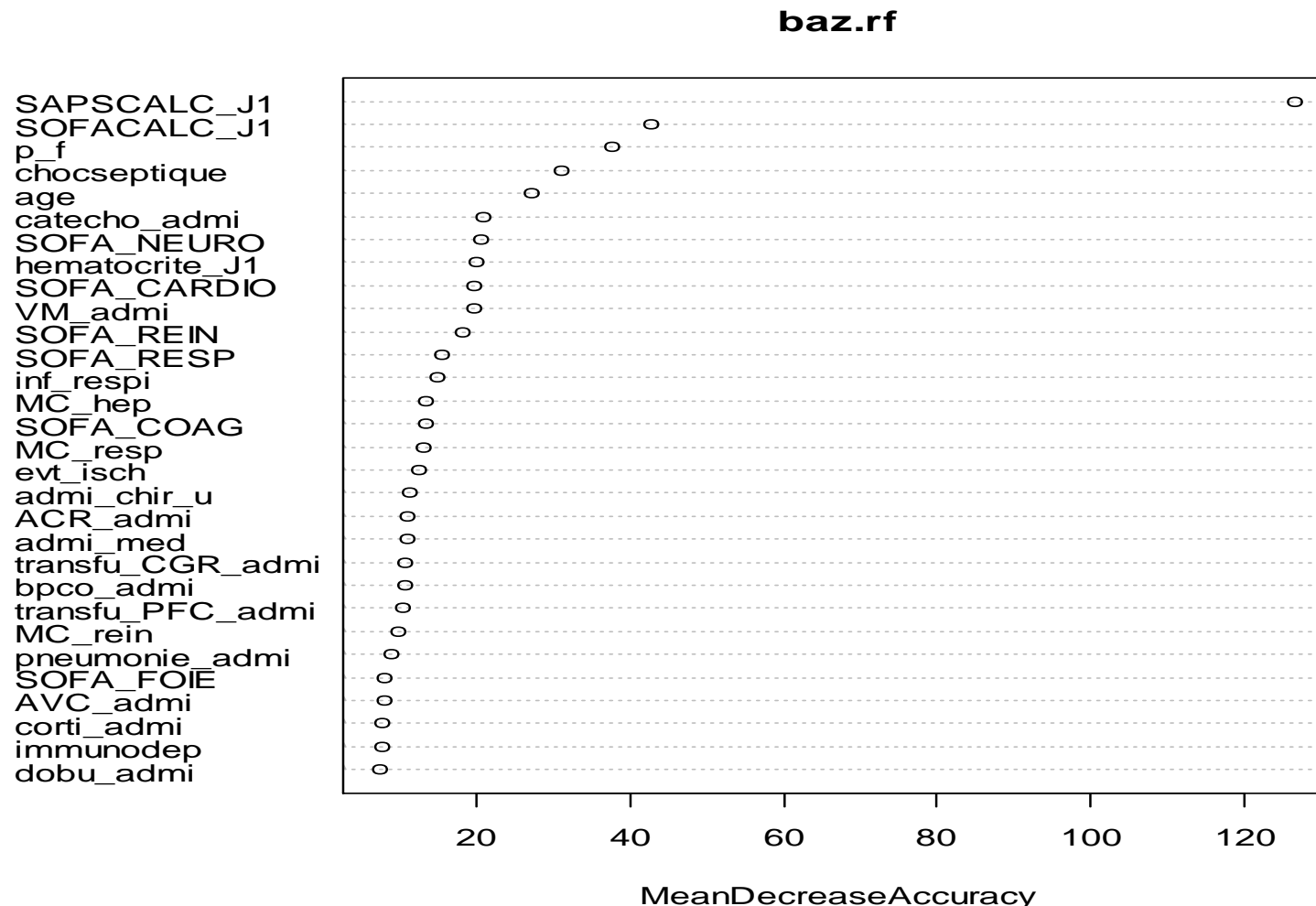
La forêt aléatoire vote pour classer un nouveau patient malade ou non malade



Apprentissage supervisé : Forêts aléatoires – Random forest



Apprentissage supervisé : Forêts aléatoires – Random forest



Décès à J30 chez les patients admis en réanimation en choc septique ou sepsis sévère en fonction des variables à l'admission – Base OUTCOMEREA – Dr Claire Dupuis

Apprentissage supervisé : Forêts aléatoires – Random forest

An Empirical Comparison of Supervised Learning Algorithms

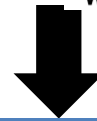
Rich Caruana
Alexandru Niculescu-Mizil
Department of Computer Science, Cornell University, Ithaca, NY 14853 USA

CARUANA@CS.CORNELL.EDU
ALEXN@CS.CORNELL.EDU

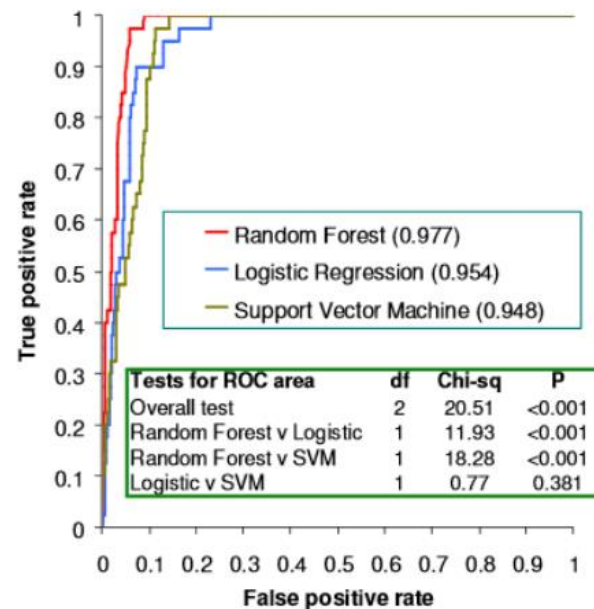
Published in final edited form as:
Int J Appl Sci Technol. 2012 August ; 2(7): 268–.

A Comparison of Logistic Regression, Logic Regression, Classification Tree, and Random Forests to Identify Effective Gene-Gene and Gene-Environmental Interactions

Wonsuk Yoo,



En faveur des forêts
aléatoires



En faveur des forêts
aléatoires

MAIS



Problème en **cas de nombre de variables pertinentes très faibles**, dans l'absolu et relativement au nombre total de variables

Manque de lisibilité du méta-modèle

Complexité de la mise en oeuvre

Introduction

Apprentissage non supervisé

Apprentissage supervisé

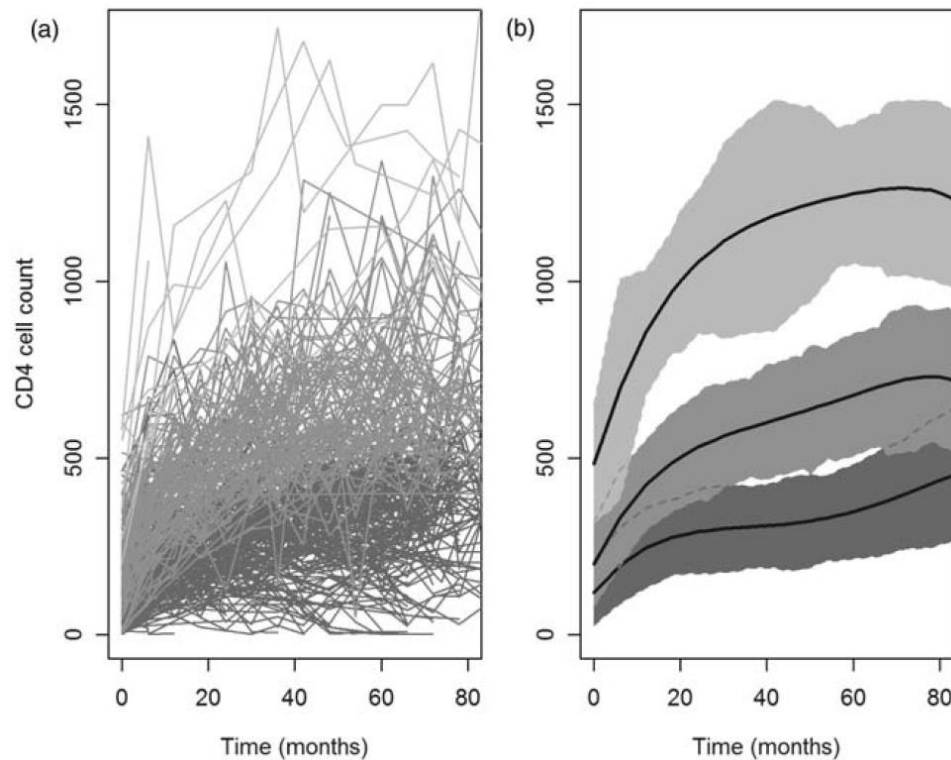


Données longitudinales

Conclusion

Group-based trajectory modeling / Latent Growth Modeling

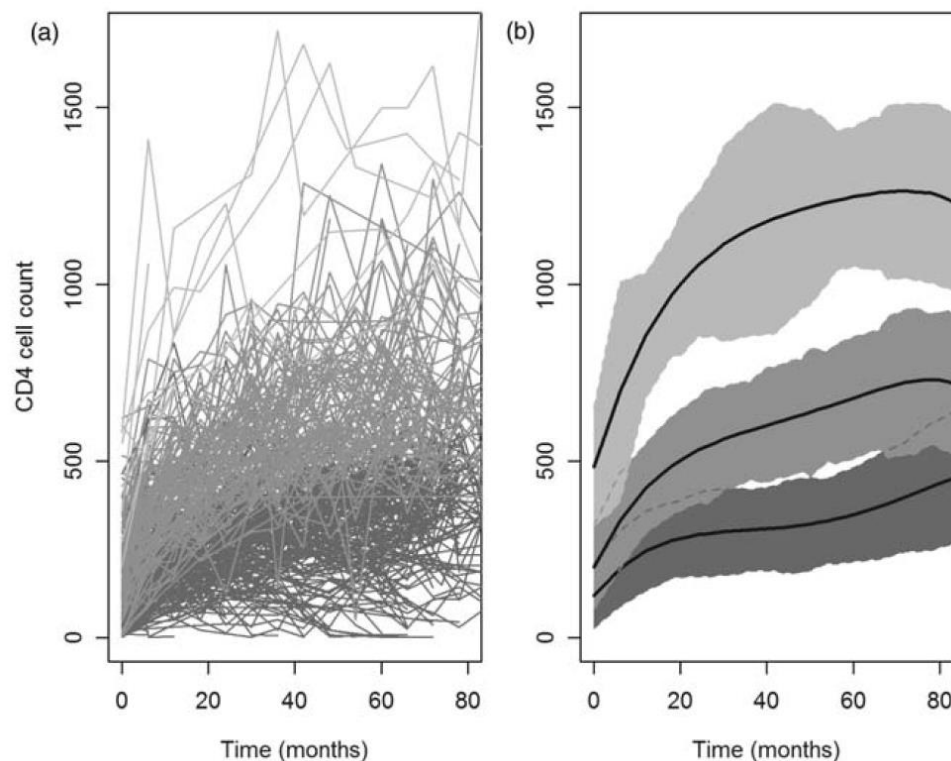
les patients sont regroupés selon leur probabilité d'appartenir à une trajectoire moyenne



Comptage des lymphocytes CD4 chez des patients VIH sous traitement antirétroviral

Group-based trajectory modeling / Latent Growth Modeling

les patients sont regroupés selon leur probabilité d'appartenir à une trajectoire moyenne



Exploitation
des **données**
longitudinales

Comptage des lymphocytes CD4 chez des patients VIH sous traitement antirétroviral

Introduction

Apprentissage non supervisé

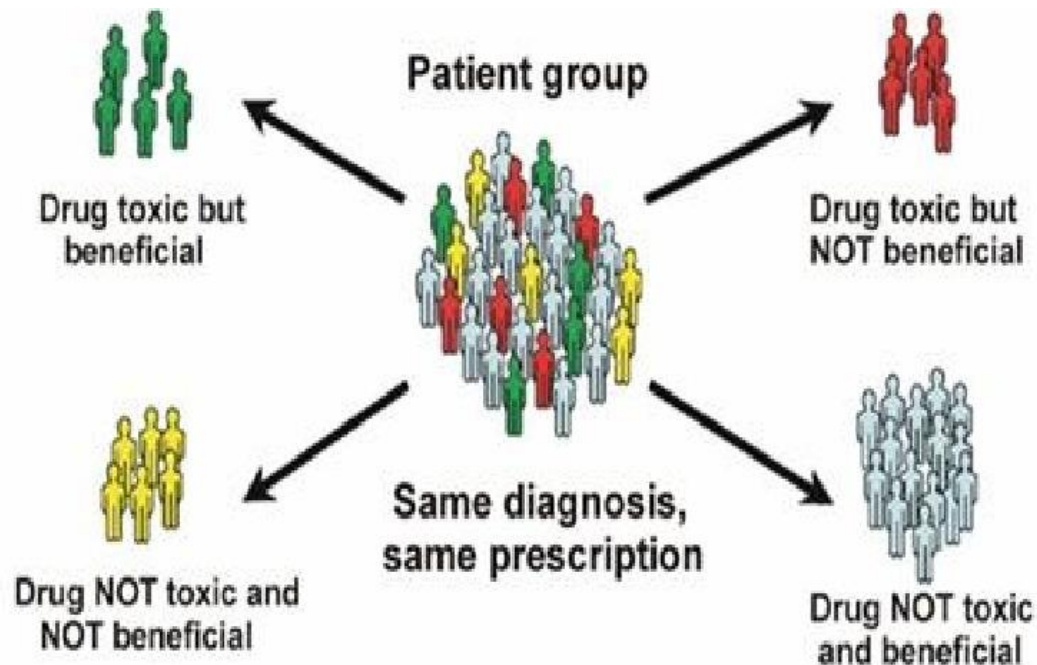
Apprentissage supervisé

Données longitudinales



Conclusion

Détermination de sous-groupes homogènes Identification de phénotypes de patients



OUTCOMEREA



Merci pour votre attention