

Apprentissage automatique

Clustering

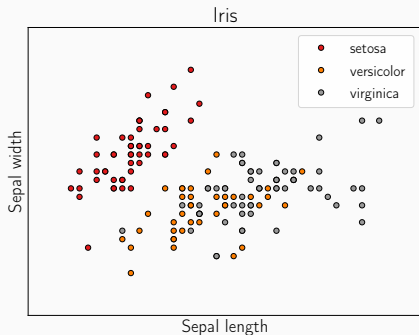
Brice Chardin

2022–2023

ISAE-ENSMA

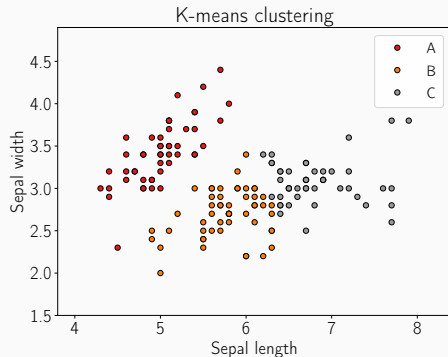
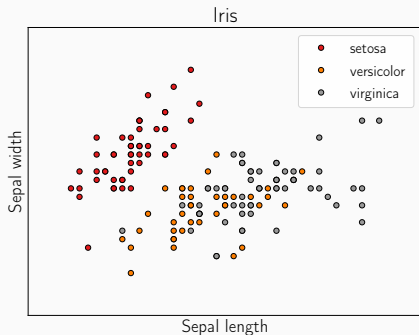
CLUSTERING

Catégoriser un jeu de données sans connaître la vérité de terrain
(i.e. apprentissage non supervisé)



CLUSTERING

Catégoriser un jeu de données sans connaître la vérité de terrain
(i.e. apprentissage non supervisé)



Partition : ensemble d'ensembles $\{P_1, \dots, P_k\}$ non nuls, disjoints et dont l'union est la population X

- $\forall i \in \{1, \dots, k\}, P_i \neq \emptyset$
- $\forall i, j \in \{1, \dots, k\}, i \neq j \implies P_i \cap P_j = \emptyset$
- $\bigcup_{i=1}^k P_i = X$

Clustering : partitionnement de la population en ensembles *homogènes* et *bien séparés*
Homogénéité et séparation dépendent des besoins

Objectif (pour certains algorithmes) : critère quantifiable à minimiser ou maximiser

Clustering sous contraintes, paramètres d'entrée du modèle

Contraintes sur les partitions :

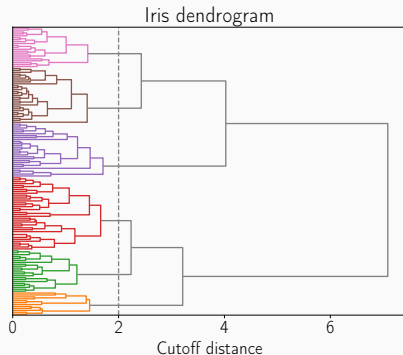
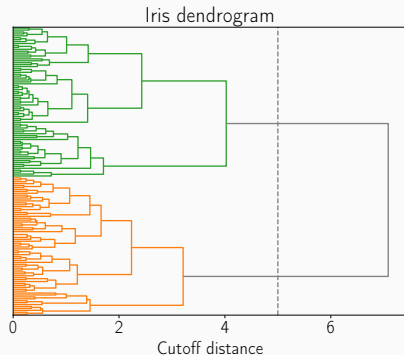
- nombre maximum de clusters
- dimensions maximales (e.g. diamètre)
- nombre minimal ou maximal d'éléments par cluster
- etc.

Contraintes sur les instances (apprentissage semi-supervisé)

- paires d'instances incompatibles (*cannot-link*)
- paires d'instances similaires (*must-link*)

CLASSIFICATION DES ALGORITHMES (1)

- Clustering hiérarchique → hiérarchie de partitions
- Clustering de partitionnement → partition unique en résultat



CLASSIFICATION DES ALGORITHMES (2)

- Méthodes de division → état initial : un seul cluster contenant la population
- Méthodes d'agglomération → état initial : un cluster par élément

Critère applicable à toute méthode hiérarchique et à certaines méthodes de partitionnement

CLASSIFICATION DES ALGORITHMES (3)

- Hard clustering → partitions discrètes (i.e. labels)
- Fuzzy clustering → degré (ou probabilité) d'appartenance à une partition

CLASSIFICATION DES ALGORITHMES (3)

- Hard clustering → partitions discrètes (i.e. labels)
- Fuzzy clustering → degré (ou probabilité) d'appartenance à une partition
- Déterministe → mêmes entrées = mêmes sorties
- Stochastique → intègre un processus aléatoire

CLASSIFICATION DES ALGORITHMES (3)

- Hard clustering → partitions discrètes (i.e. labels)
- Fuzzy clustering → degré (ou probabilité) d'appartenance à une partition
- Déterministe → mêmes entrées = mêmes sorties
- Stochastique → intègre un processus aléatoire

Si objectif défini :

- Exact → garantit une solution optimale
- Approximatif → ne garantit pas une solution optimale

CLUSTERING EN K-MOYENNES (*K-MEANS*)

Contrainte : nombre prédéfini k de clusters

Objectif : minimiser l'écart quadratique moyen du clustering

Soient $\{C_1, \dots, C_k\}$ clusters de centres $\{\mu_1, \dots, \mu_k\}$

$$\text{écart} = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

CLUSTERING EN K-MOYENNES (*K-MEANS*)

Contrainte : nombre prédéfini k de clusters

Objectif : minimiser l'écart quadratique moyen du clustering

Soient $\{C_1, \dots, C_k\}$ clusters de centres $\{\mu_1, \dots, \mu_k\}$

$$\text{écart} = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

Solution optimale difficile à calculer

ALGORITHME DES K-MOYENNES (*K-MEANS*)

Algorithme itératif simple :

1. affecter un point à chaque cluster
2. répéter :
 - calculer le centre de chaque cluster
 - affecter chaque point au cluster dont le centre est le plus proche
3. jusqu'à convergence (i.e. pas de réaffectation)

ALGORITHME DES K-MOYENNES (*K-MEANS*)

Algorithme itératif simple :

1. affecter un point à chaque cluster
2. répéter :
 - calculer le centre de chaque cluster
 - affecter chaque point au cluster dont le centre est le plus proche
3. jusqu'à convergence (i.e. pas de réaffectation)

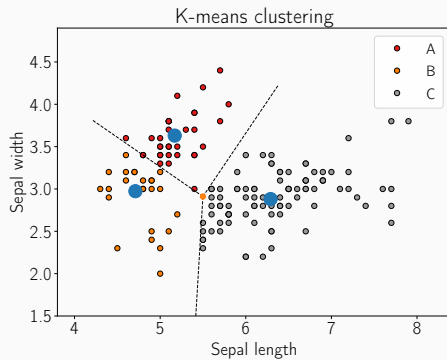
Convergence garantie, mais vers un minimum local (solution approximative)

Résultat (très) dépendant de l'initialisation

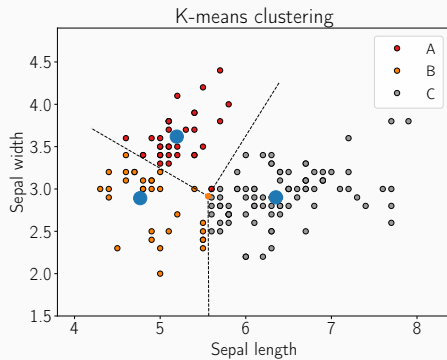
- meilleur résultat de plusieurs exécutions avec initialisations différentes
- heuristiques d'initialisation (e.g. k-means++)

Parmi les algorithmes de partitionnement les plus rapides

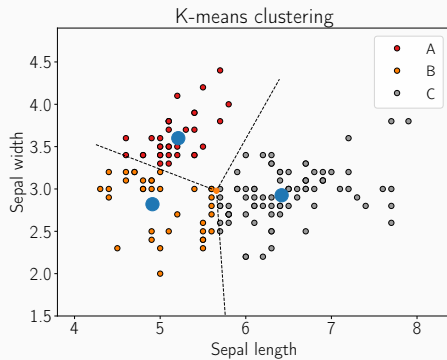
ITÉRATIONS DE K-MEANS



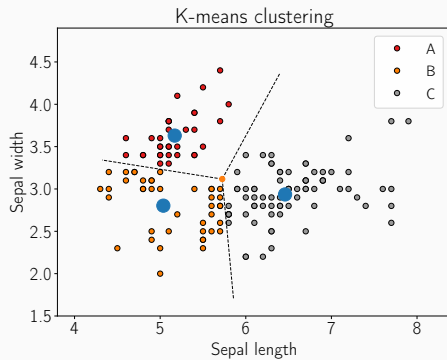
ITÉRATIONS DE K-MEANS



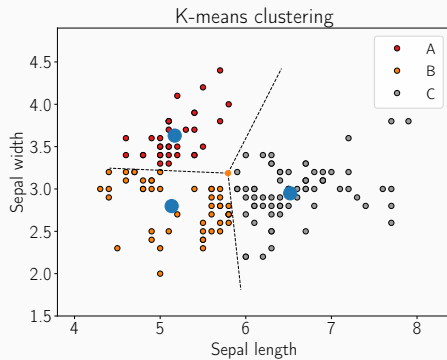
ITÉRATIONS DE K-MEANS



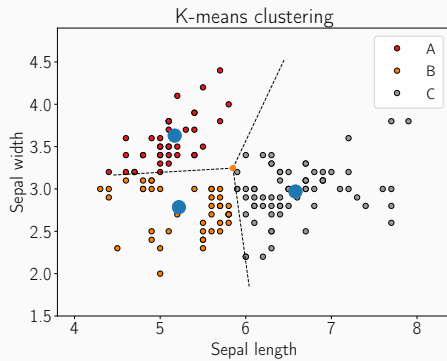
ITÉRATIONS DE K-MEANS



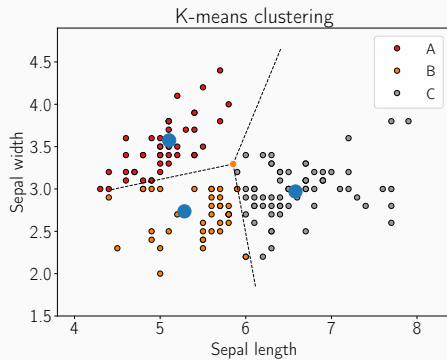
ITÉRATIONS DE K-MEANS



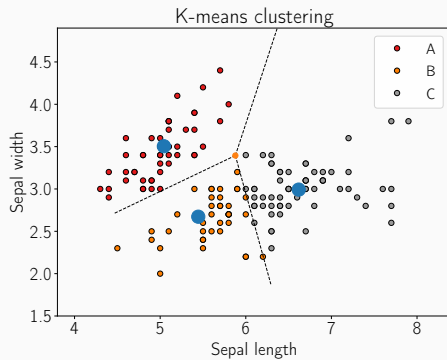
ITÉRATIONS DE K-MEANS



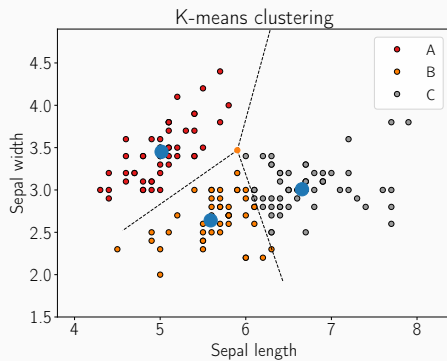
ITÉRATIONS DE K-MEANS



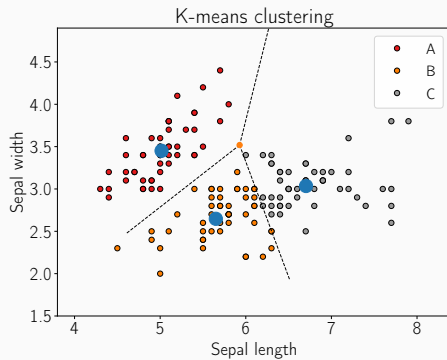
ITÉRATIONS DE K-MEANS



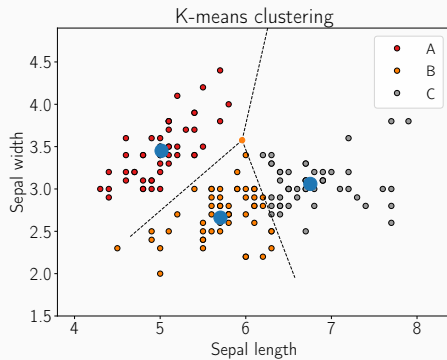
ITÉRATIONS DE K-MEANS



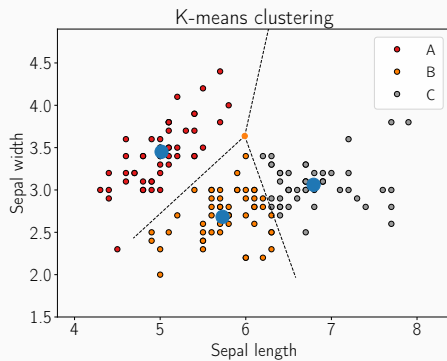
ITÉRATIONS DE K-MEANS



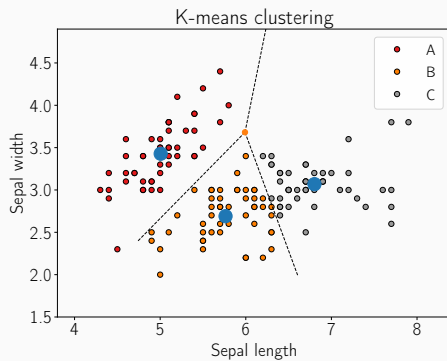
ITÉRATIONS DE K-MEANS



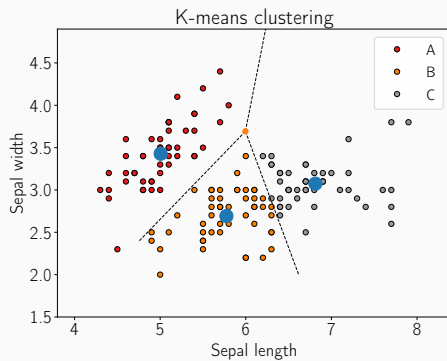
ITÉRATIONS DE K-MEANS



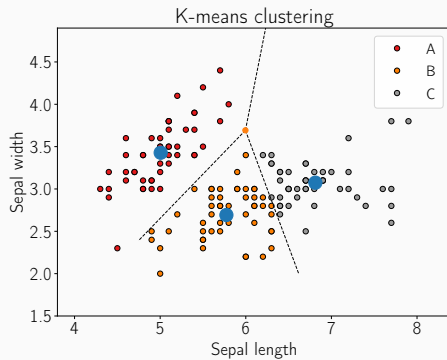
ITÉRATIONS DE K-MEANS



ITÉRATIONS DE K-MEANS



ITÉRATIONS DE K-MEANS



→ Convergence

Même principe que k-means en remplaçant les centres par des médoïdes

Utile lorsqu'une moyenne n'est pas calculable

Données catégorielles

- $d(\text{pomme}, \text{pomme}) = 0$
- $d(\text{pomme}, \text{poire}) = 1$
- $\text{avg}(\text{pomme}, \text{poire}) = ?$

Définition (médoïde)

$$\text{médoïde}(C) = \underset{c \in C}{\operatorname{argmin}} \sum_{e \in C} d(e, c)$$

I.e. point du cluster dont la distance moyenne aux autres membres est minimale

Exemple d'algorithme approximatif : PAM (*Partitioning Around Medoids*)

MÉTHODES PAR DENSITÉ : DBSCAN ET OPTICS

Méthodes par densité, définies par deux paramètres : ϵ et *minpts*

Zone dense : zone contenant au moins *minpts* points dans un rayon de ϵ

Propagation de la notion de densité pour obtenir les clusters

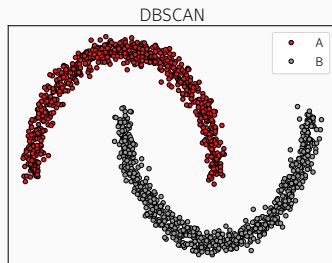
MÉTHODES PAR DENSITÉ : DBSCAN ET OPTICS

Méthodes par densité, définies par deux paramètres : ϵ et *minpts*

Zone dense : zone contenant au moins *minpts* points dans un rayon de ϵ

Propagation de la notion de densité pour obtenir les clusters

Permet d'obtenir des formes non convexes



CLASSIFICATION ASCENDANTE HIÉRARCHIQUE

1. affecter chaque point à un cluster distinct
2. tant qu'il existe plusieurs clusters
 - fusionner les deux clusters les plus *proches*

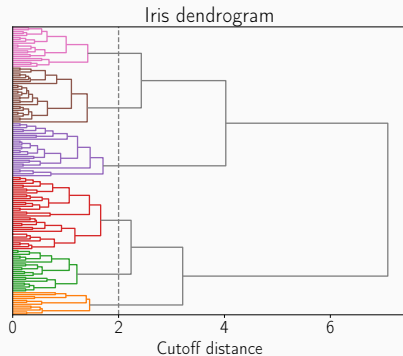
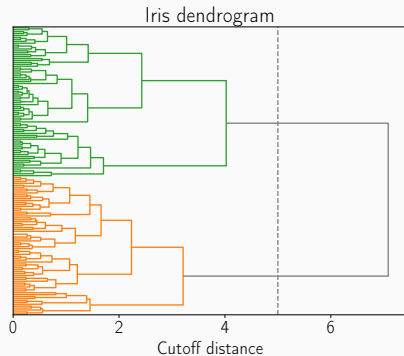
Critères de proximité entre deux clusters A et B :

- distance maximale entre deux points (complete-link) : $\max\{d(a, b) : a \in A, b \in B\}$
- distance minimale entre deux points (single-link) : $\min\{d(a, b) : a \in A, b \in B\}$
- distance entre les médoïdes de A et B

- distance moyenne :
$$\frac{\sum_{a \in A} \sum_{b \in B} d(a, b)}{|A| \cdot |B|}$$
- etc.

CLASSIFICATION ASCENDANTE HIÉRARCHIQUE

Critère de coupure sur la hiérarchie résultat : distance ou nombre de clusters



MESURES DE QUALITÉ

TAUX DE BONNE CLASSIFICATION

Si vérité de terrain connue : $T_x = \frac{\text{nb éléments bien classés}}{\text{nb éléments total}}$

Matrice de contingence

Classes	Clusters			
	C ₁	C ₂	C ₃	C ₄
A	15	7	10	0
B	12	2	5	0
C	3	20	3	11
D	0	4	12	10

TAUX DE BONNE CLASSIFICATION

Si vérité de terrain connue : $T_x = \frac{\text{nb éléments bien classés}}{\text{nb éléments total}}$

Matrice de contingence

Classes	Clusters			
	C ₁	C ₂	C ₃	C ₄
A	15	7	10	0
B	12	2	5	0
C	3	20	3	11
D	0	4	12	10

Considérer le cas le plus favorable

TAUX DE BONNE CLASSIFICATION

Si vérité de terrain connue : $T_x = \frac{\text{nb éléments bien classés}}{\text{nb éléments total}}$

Matrice de contingence

Classes	Clusters			
	C ₁	C ₂	C ₃	C ₄
A	15	7	10	0
B	12	2	5	0
C	3	20	3	11
D	0	4	12	10

Considérer le cas le plus favorable

$$T_x = \frac{10 + 12 + 20 + 10}{114} = 45.6\%$$

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

Classes	Clusters
A	C ₁
B	C ₂
A	C ₂
B	C ₂
C	C ₁

Concordances : 0

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

	Classes	Clusters	
\neq	A	C ₁	\neq
	B	C ₂	
	A	C ₂	
	B	C ₂	
	C	C ₁	

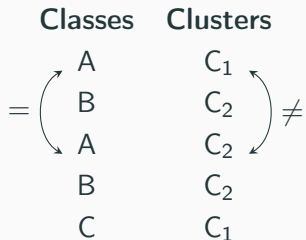
Concordances : 1

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :



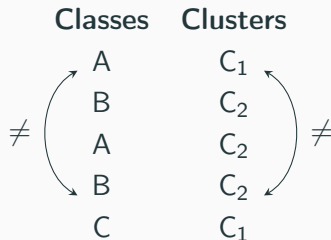
Concordances : 1

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :



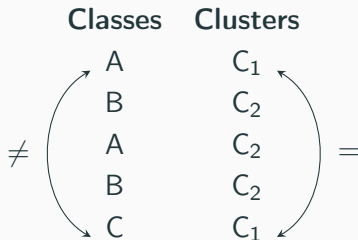
Concordances : 2

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :



Concordances : 2

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

	Classes	Clusters	
	A	C ₁	
	B	C ₂	
≠	A	C ₂	↗ ↘ =
	B	C ₂	
	C	C ₁	

Concordances : 2

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

	Classes	Clusters	
	A	C ₁	
	B	C ₂	
=	A	C ₂	=
	B	C ₂	
	C	C ₁	

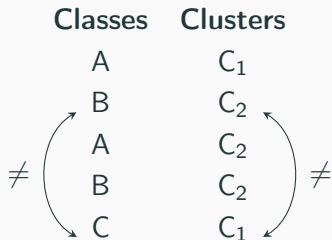
Concordances : 3

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :



Concordances : 4

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

	Classes	Clusters	
	A	C ₁	
	B	C ₂	
≠	A	C ₂	⌋ =
	B	C ₂	
	C	C ₁	

Concordances : 4

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

	Classes	Clusters	
	A	C ₁	
	B	C ₂	
\neq	A	C ₂	\neq
	B	C ₂	
	C	C ₁	

Concordances : 5

Total : 10

INDICE DE RAND

Proportion de paires d'éléments égaux ou différents dans les deux cas

Indice de Rand non ajusté :

Classes	Clusters
A	C ₁
B	C ₂
A	C ₂
B	C ₂
≠ (C	C ₁) ≠

Concordances : 6

Total : 10

Indice de Rand : $\frac{6}{10} = 0.6$

Indice de Rand non ajusté

- Score parfait de 1 (concordance totale)
- Pire score de 0 (discordance totale) difficile à obtenir

Indice de Rand non ajusté

- Score parfait de 1 (concordance totale)
- Pire score de 0 (discordance totale) difficile à obtenir

→ Les concordances sur les différences donnent des scores proches de 1 même pour des clusters très dissimilaires

Indice de Rand non ajusté

- Score parfait de 1 (concordance totale)
- Pire score de 0 (discordance totale) difficile à obtenir

→ Les concordances sur les différences donnent des scores proches de 1 même pour des clusters très dissimilaires

Indice de Rand ajusté : score entre -1 et 1, affectation aléatoire → 0

Autres mesure si la vérité de terrain n'est pas connue

Inertie intra-cluster (mesure d'homogénéité des clusters)

$$I_{\text{intra}} = \frac{1}{|X|} \sum_{i=1}^k \sum_{x \in C_i} d^2(x, \mu_i) \quad \leftarrow \text{fonction objectif de } k\text{-means}$$

Score de silhouette : $s(x_i) = \frac{b(x_i) - a(x_i)}{\max(a(x_i), b(x_i))}$

$a(x_i)$: distance moyenne entre $x_i \in C_m$ et les points appartenant au même cluster

$b(x_i)$: minimum de la distance moyenne entre x_i et les points d'un autre cluster

SCORE DE SILHOUETTE

Score de silhouette : $s(x_i) = \frac{b(x_i) - a(x_i)}{\max(a(x_i), b(x_i))}$

$$a(x_i) = \frac{1}{|C_m| - 1} \sum_{x_j \in C_m, j \neq i} d(x_i, x_j) \qquad b(x_i) = \min_{l \neq m} \left(\frac{1}{|C_l|} \sum_{x_j \in C_l} d(x_i, x_j) \right)$$

Si $s(x_i) < 0$, alors $b(x_i) < a(x_i)$: x_i pourrait (ou devrait) changer de cluster

Score de silhouette global correspond au score de silhouette moyen

$$s(X) = \frac{1}{|X|} \sum_{x_i \in X} s(x_i)$$

Un score de silhouette parfait (i.e. $a(x_i) = 0$, généralement inatteignable) vaut 1