

Analyse de Données – Développeur

C5-160512-INFO

4. Classification Non Supervisée

Pr. Carl FRÉLICOT

FST/PAS221 – carl.frelicot@univ-lr.fr



Licence d'Informatique 3ème année

Automne 2021 - La Rochelle



1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation



- individus vs variables quantitatives
données dites *vectorielles*

$$X = [X_{ki}]_{k=1,n; i=1,p}$$

nuage de n points en dim. p

données sans étiquette de groupe

cas non supervisé

- individus vs individus
données non vectorielles

ex. : tableau de distances, proximités, etc

$$D = [d_{kl}]_{k,l=1,n} \text{ où } d_{kl} = d(x_k, x_l)$$

K-Means, Hierarchical Clustering

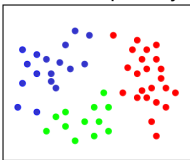
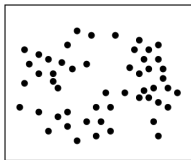
- variables vs variables

$$D = [d_{ij}]_{i,j=1,p} \text{ où } d_{ij} = d(X_i, X_j)$$

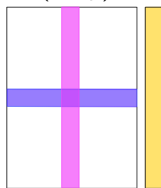
Feature Clustering

par exemple ?

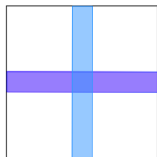
But : identifier une structure de groupes dans X
afin d'attribuer à tout x une étiquette y



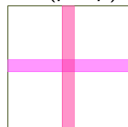
$X (n \times p)$ $Y (n \times 1)$



$D (n \times n)$



$D (p \times p)$





- différents types de classification non supervisée selon qu'on cherche :
 - à modéliser les classes par des distributions conditionnellement aux groupes (approche probabiliste)
 - une structure particulière (approche métrique)
 - ▶ partition
 - ▶ hiérarchie
 - ▶ recouvrement

*mixture
decomposition*

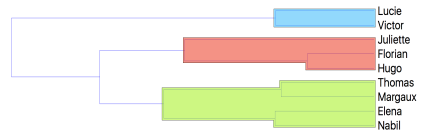
*K-Means
Hierarchical Clustering
plusieurs étiquettes*

Exemple : notes de $n = 9$ élèves de Terminale S dans $p = 5$ matières

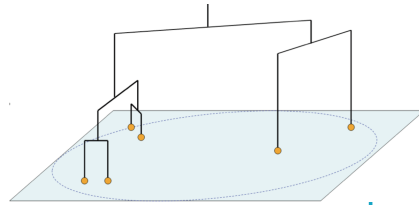
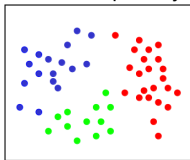
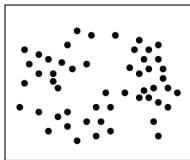
	Prénom	math.	info.	fran.	angl.	arts
1	Thomas	6.0	6.0	5.0	5.5	8.0
2	Margaux	8.0	8.0	8.0	8.0	9.0
3	Florian	6.0	7.0	11.0	10.0	11.0
4	Lucie	15.0	14.0	16.0	15.0	8.0
5	Victor	14.0	14.0	12.0	12.5	10.0
6	Elena	11.0	10.0	5.5	7.0	13.0
7	Hugo	5.5	7.0	14.0	11.5	10.0
8	Nabil	13.0	12.5	8.5	9.5	12.0
9	Juliette	9.0	9.5	12.5	12.0	18.0

	Prénom	Cluster
1	Thomas	C3
2	Margaux	C3
3	Florian	C2
4	Lucie	C1
5	Victor	C1
6	Elena	C3
7	Hugo	C2
8	Nabil	C3
9	Juliette	C2

lien entre partition et hiérarchie ?



But : identifier une structure de groupes dans X afin d'attribuer à tout x une étiquette y





1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation



Distances entre individus : une **distance métrique**^(a) d sur un ensemble X :

$X \times X \rightarrow \mathbb{R}_+, (x, y) \mapsto d(x, x)$ vérifiant

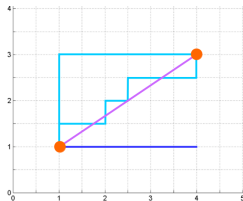
- (1) $d(x, y) = d(y, x)$ *symétrie*
- (2) $d(x, y) = 0 \Leftrightarrow x = y$ *indiscernabilité*
- (3) $d(x, z) \leq d(x, y) + d(y, z)$ *inégalité triang.*

^(a) si tout $x \in X$ se projette dans un espace (vectoriel) \mathcal{X} , (\mathcal{X}, d) est appelé *espace métrique* ; pour nous $\mathcal{X} = \mathbb{R}^p$, d'où le raccourci de langage entre **distance** et **métrique**

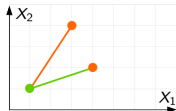
- une **pseudométrique** ne vérifie que (1)-(3)
- une **semimétrique** ou **dissemblance** ne vérifie que (1)-(2) *suffit en Classif. NS*
- une **ultramétrique** vérifie (1)-(2)-(3') $d(x, z) \leq \max(d(x, y), d(y, z))$ *distance ?*

Distances usuelles : distances de **Minkowski** $d_q(x, y) = \left(\sum_{j=1}^p |x_j - y_j|^q \right)^{1/q}$

- *Manhattan* ou *cityblock* si $q = 1$
- *euclidienne* si $q = 2$; alors $d_2^2(x, y) = \|x - y\|_2^2$
...
- *Chebychev* si $q \rightarrow +\infty$;
alors $d_\infty(x, y) = \max_{j=1, p} |x_j - y_j|$



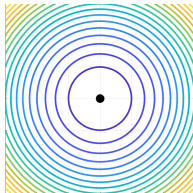
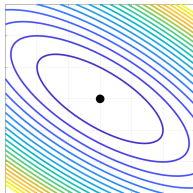
Pondération des variables !





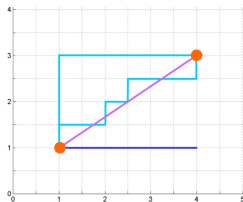
Distance de Mahalanobis : $d_M^2(x, y) = \|x - y\|_{V^{-1}}^2 = {}^t(x - y) V^{-1} (x - y)$

- si V est diagonale, alors les ellipses d'iso-distance sont parallèles aux axes
- si $V = I$? alors $d_M^2(x, y) = d_2^2(x, y)$

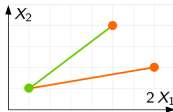
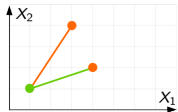


Distances usuelles : distances de Minkowski $d_q(x, y) = \left(\sum_{j=1}^p |x_j - y_j|^q \right)^{1/q}$

- Manhattan ou cityblock si $q = 1$
- euclidienne si $q = 2$; alors $d_2^2(x, y) = \|x - y\|_2^2$
- Chebychev si $q \rightarrow +\infty$;
alors $d_\infty(x, y) = \max_{j=1,p} |x_j - y_j|$



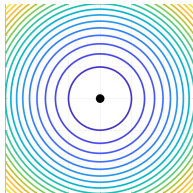
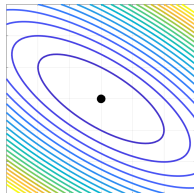
Pondération des variables !





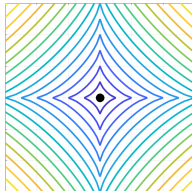
Distance de Mahalanobis : $d_M^2(x, y) = \|x - y\|_{V^{-1}}^2 = {}^t(x - y) V^{-1} (x - y)$

- si V est diagonale, alors les ellipses d'iso-distance sont parallèles aux axes
- si $V = I$? alors $d_M^2(x, y) = d_2^2(x, y)$

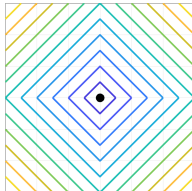


$q = 2$

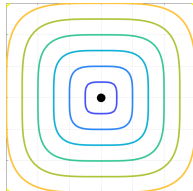
Distances usuelles : distances de Minkowski $d_q(x, y) = \left(\sum_{j=1}^p |x_j - y_j|^q \right)^{1/q}$



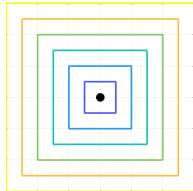
$q ?$ $q = 0.75$



$q = 1$



$q = 4$



$q = +\infty$



1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.3.1 Motivation
 - 4.3.2 Algorithme C-means
 - 4.3.3 Pros and Cons – C-means
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation



Partition : une **partition stricte** d'un ensemble $X = \{x_1, x_2, \dots, x_n\}$ est un ensemble $P = \{C_1, C_2, \dots, C_c\}$ de groupes (*clusters*) disjoints, couvrant X

- on lui associe une **matrice de partition** $U = [u_{ik}]_{i=1,c;k=1,n}$ où
 où $u_{ik} = 1$ ou 0 selon que $x_k \in C_i$ ou non,
 telle que $\sum_i u_{ik} = 1, \forall k$
 et $0 < \sum_k u_{ik} < n, \forall i$

stricte

Problèmes :

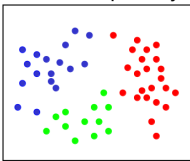
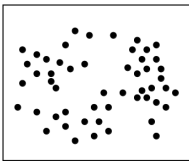
- nombre c de groupes ?
- la recherche de la "meilleure partition" en k classes d'un ensemble de n éléments est NP-complet ;
 nombre de Stirling :

$$S(c, n) = \frac{1}{c!} \sum_{i=1}^c (-1)^{c-i} \binom{i}{c} i^n \quad S(c, n) \simeq \frac{c^n}{c!}$$

c et n grands

c	n	S
2	4	7
2	10	511
5	10	42 535
2	30	536 870 911
3	150	$> 6 \times 10^{70}$

But : identifier une structure de groupes dans X
 afin d'attribuer à tout x une étiquette y

 $X \rightarrow U$ ou P 

Sont nécessaires un(e) :

- 1 dissimilarité/distance d sur X
- 2 critère d'affectation \mathcal{D} aux groupes
- 3 algorithme d'optimisation de \mathcal{D}



1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.3.1 Motivation
 - 4.3.2 Algorithme C-means
 - 4.3.3 Pros and Cons – C-means
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation



Algorithme(s) C-means

- méthode de **clustering** : recherche itérative d'une partition

$$X \rightarrow U, V$$

et de représentants (?) $V = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_c\}$

- paramètres utilisateurs :

- nombre c de clusters
- paramètre d'arrêt ε
- ① distance d

attention : $K == c$

ou nombre max. d'itérations

d_2 par défaut

- critère $\mathcal{D}(U, V) = \frac{1}{n} \sum_{i=1}^c \sum_{x_k \in C_i} d_2^2(x_k, \bar{x}_i)$

inertie intra-clusters

min ? max ?

ou $V^{(0)}$

- ① initialisation aléatoire : $U^{(0)}$

stabilité de la partition

- ① tant que non convergence

① $V^{(t)} = \operatorname{argmin}_V \mathcal{D}(U^{(t-1)}, V)$

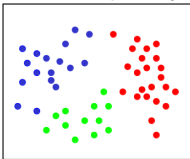
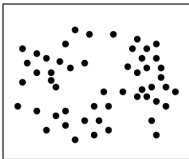
prototypage

② $U^{(t)} = \operatorname{argmin}_U \mathcal{D}(U, V^{(t)})$

affectation

But : identifier une structure de groupes dans X
afin d'attribuer à tout x une étiquette y

$$X \rightarrow U \text{ ou } P$$



Sont nécessaires un(e) :

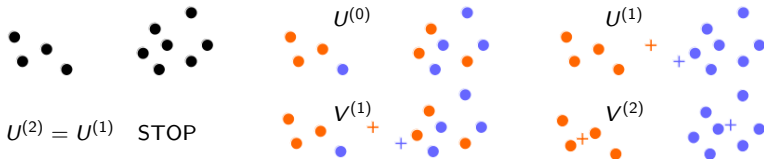
- ① dissimilarité/distance d sur X
- ② critère d'affectation \mathcal{D} aux groupes
- ③ algorithme d'optimisation de \mathcal{D}



Algorithm(e)s C-means

K-moyennes ou "centres mobiles"

- méthode de **clustering** : recherche itérative d'une partition $X \rightarrow U, V$
et de représentants (?) $V = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_c\}$
- paramètres utilisateurs :
 - nombre c de clusters
 - paramètre d'arrêt ε
 - ① distance d
- critère $\mathcal{D}(U, V) = \frac{1}{n} \sum_{i=1}^c \sum_{x_k \in C_i} d_2^2(x_k, \bar{x}_i) = \frac{1}{n} \sum_{k=1}^n d_2^2(x_k, \bar{x}_{C_i(x_k)})$ inertie intra
min \rightarrow max \rightarrow
ou $V^{(0)}$
- ① initialisation aléatoire : $U^{(0)}$
- ① tant que non convergence
 - ① $V^{(t)} = \operatorname{argmin}_V \mathcal{D}(U^{(t-1)}, V)$ prototypage
 - ② $U^{(t)} = \operatorname{argmin}_U \mathcal{D}(U, V^{(t)})$ affectation
- convergence : $\|U^{(t)} - U^{(t-1)}\| < \varepsilon$ norme matricielle, par ex. : \max ou \maxiter
- $U^{(0)} \rightarrow V^{(1)} \rightarrow U^{(1)} \rightarrow V^{(2)} \dots V^{(t)} = V^{(t-1)}$
 $V^{(0)} \rightarrow U^{(1)} \rightarrow V^{(1)} \rightarrow U^{(2)} \dots U^{(t)} = U^{(t-1)}$

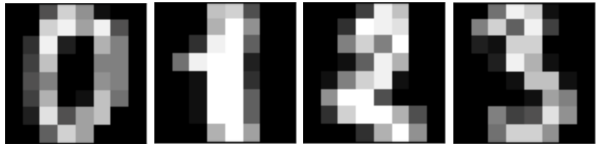




Exemple : chiffres manuscrits

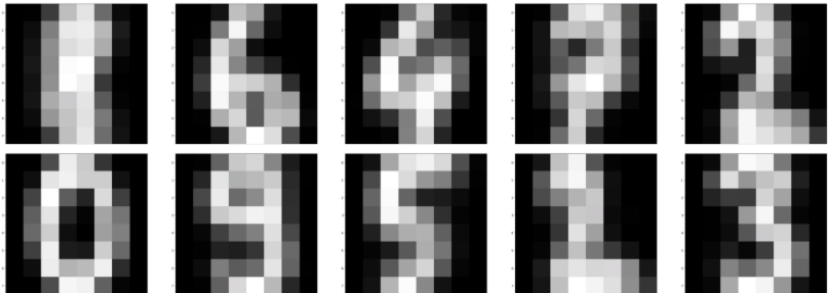
ex. images originales

- 10 chiffres
1 797 imageries
($\simeq 180$ par chiffre)
de très petite taille
 $64 = 8 \times 8$ pixels



- barycentres trouvés

8, 6, 4, 7, 2



- taux d'accord avec la **vérité-terrain** : 92%
supervision connue, par expertise / connaissance humaine

0, 9, 5, 1, 3
ground-truth



1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.3.1 Motivation
 - 4.3.2 Algorithme C-means
 - 4.3.3 Pros and Cons – C-means
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation

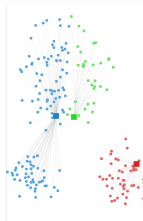


- ⊕ intuitif et rapide
- ⊖ sensible à l'initialisation
solution : exécuter r fois,
retenir les points qui finissent toujours ensemble

comparé à *HAC*
optimum local

formes fortes

run #1

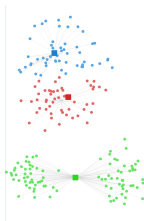
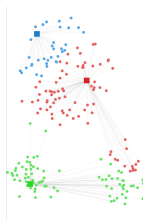


initialisation



partition finale

run #2



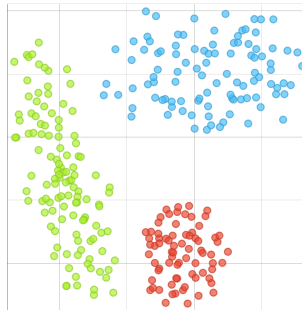
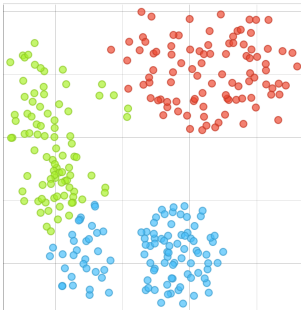


- ⊕ intuitif et rapide
- ⊖ sensible à l'initialisation
solution : exécuter r fois,
retenir les points qui finissent toujours ensemble
- ⊖ tendance à former des clusters équilibrés
 - effectif
 - dispersionsolution : adapter la métrique, donc le critère
- ⊖ tendance à former des clusters hypersphériques
solution : adapter la métrique

comparé à *HAC*
optimum local

formes fortes
inertie intra-classes

variantes
distance euclidienne
une idée ?





4.3 Méth. par Partition

4.3.3 Pros and Cons

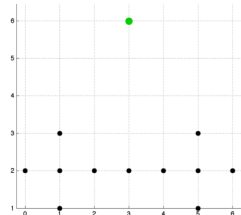
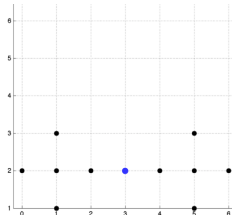
- ⊕ intuitif et rapide
- ⊖ sensible à l'initialisation
solution : exécuter r fois,
retenir les points qui finissent toujours ensemble
- ⊖ tendance à former des clusters équilibrés
 - effectif
 - dispersionsolution : adapter la métrique, donc le critère
- ⊖ tendance à former des clusters hypersphériques
solution : adapter la métrique
- ⊖ sensible aux *inliers*
et *outliers*

comparé à *HAC*
optimum local

formes fortes
inertie intra-classes

variantes
distance euclidienne
une idée ?

partition stricte
 C fixé : *cluster validity*





Cluster Validity Problem

- quel c_{opt} ?
- critères (non/liés à l'algorithme)
fondés sur les points et/ou centres

$$c_{opt.} = \operatorname{argmin}_{c=2, c_{max}} \operatorname{max}_{c=2, c_{max}} CVI(c)$$

Dunn Index :

$$DI(c) = \frac{\min_{1 \leq i < i' \leq c} d(\bar{x}_i, \bar{x}_{i'})}{\max_{j=1, c} \Delta_j} \quad \text{où } \Delta_j = \max_{x_k, x_l \in C_j} d(x_k, x_l)$$

certains CVI ont une tendance monotone avec c :(

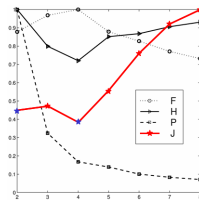
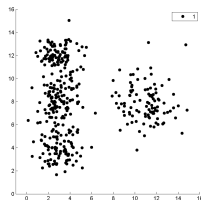
- critères relatifs (comparaison de deux partitions en c et c' clusters)

$$\text{Rand Index} : RI(P, Q) = \frac{2t - (u+v)}{n(n-1)} + 1$$

- $t = \sum_{i=1}^c \sum_{j=1}^{c'} n_{ij}^2 - n$
- $u = \sum_{i=1}^c n_{i\bullet}^2 - n$, où $n_{i\bullet} = \sum_{j=1}^{c'} n_{ij}$
- $v = \sum_{j=1}^{c'} n_{\bullet j}^2 - n$, où $n_{\bullet j} = \sum_{i=1}^c n_{ij}$

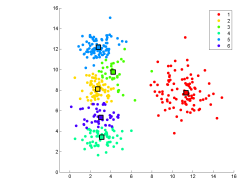
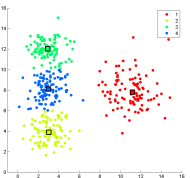
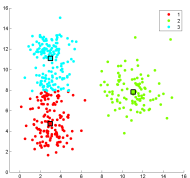
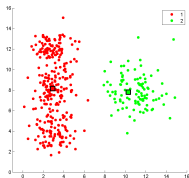
à partir de la matrice d'accord :

$$N(P, Q) = {}^t P Q$$

de dimension $(p \times q)$ 

min ou max ?

min ou max ?





1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation



- individus vs variables quantitatives
données dites *vectorielles*

$$X = [X_{ki}]_{k=1,n; i=1,p}$$

nuage de n points en dim. p

données sans étiquette de groupe

cas non supervisé

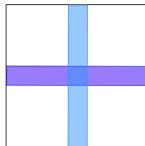
- individus vs individus
données non vectorielles

ex. : tableau de distances, proximités, etc

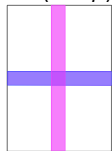
$$D = [d_{kl}]_{k,l=1,n} \text{ où } d_{kl} = d(x_k, x_l)$$

C-Means, Hierarchical Clustering

$D (n \times n)$



$X (n \times p)$



- variables vs individus
données *vectorielles*

$${}^tX = [X_{ik}]_{i=1,p; k=1,n}$$

nuage de p points en dim. n

- variables vs variables
données non vectorielles

ex. : tableau de distances, proximités, etc

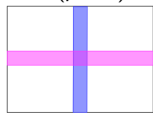
$$D = [d_{ij}]_{i,j=1,p} \text{ où } d_{ij} = d(X_i, X_j)$$

Feature Clustering (C-Means, Hierarchical)

$D (p \times p)$



${}^tX (p \times n)$



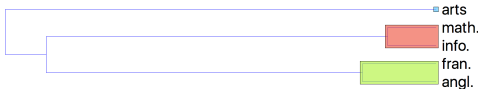
par ex. ?



Exemple : notes de $n = 9$ élèves de Terminale S dans $p = 5$ matières

Feature name	Thomas	Margaux	Florian	Lucie	Victor	Elena	Hugo	Nabil	Juliette
math.	6.000	8.000	6.000	15.000	14.000	11.000	5.500	13.000	9.000
info.	6.000	8.000	7.000	14.000	14.000	10.000	7.000	12.500	9.500
fran.	5.000	8.000	11.000	16.000	12.000	5.500	14.000	8.500	12.500
angl.	5.500	8.000	10.000	15.000	12.500	7.000	11.500	9.500	12.000
arts	8.000	9.000	11.000	8.000	10.000	13.000	10.000	12.000	18.000

- nuage de variables ($p = 5 \times n = 9$)
- algorithme *C-means* ; $c = 3$



- utiliser les centres comme variables
→ espace de dimension réduite ($q = c < p$)

- variables vs individus
données *vectorielles*

$${}^tX = [X_{ik}]_{i=1,p; k=1,n}$$

nuage de p points en dim. n

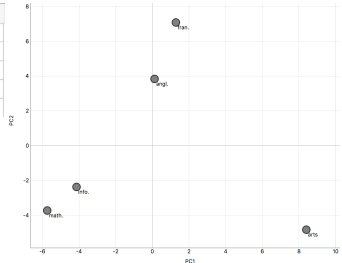
- variables vs variable
données non vectorielles

ex. : tableau de distances, proximités, etc

$$D = [d_{ij}]_{i,j=1,p} \text{ où } d_{ij} = d(X_i, X_j)$$

Feature Clustering (*C-Means*, *Hierarchical*)

par ex. ?



ACP sur les variables (plan 1-2)

$${}^tX (p \times n)$$



$$D (p \times p)$$





1. Préambule
2. Tableaux et Espaces
3. Réduction de la Dimensionnalité
4. Classification Non Supervisée
 - 4.1 Problématique et Structures
 - 4.2 Distances
 - 4.3 Méthodes de Partitionnement (Clustering)
 - 4.3.1 Motivation
 - 4.3.2 Algorithme C-means
 - 4.3.3 Pros and Cons – C-means
 - 4.4 Réduction de la Dimensionnalité ?
5. Classification Supervisée
6. Conception et Évaluation