

Cours: VISUALISATION DE DONNEES

Master MIV, 2022/2023

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

Sommaire

5.1 Les données temporelles

5.2 Clustering de données temporelles

5.3 Cas d'images de télédétection

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

Les données sont enregistrés de manière continue pour la biologie, la télédétection, la santé, etc.

Caractériser les données chronologiques pour comprendre l'évolution et les relations temporelles des variables dans le temps et prédire leur comportement futur.

L'exploitation de cette dimension temporelle introduit une complexité supplémentaire dans les différentes tâches:

- Fouille de données
- Extraction de connaissances
- **Visualisation**

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

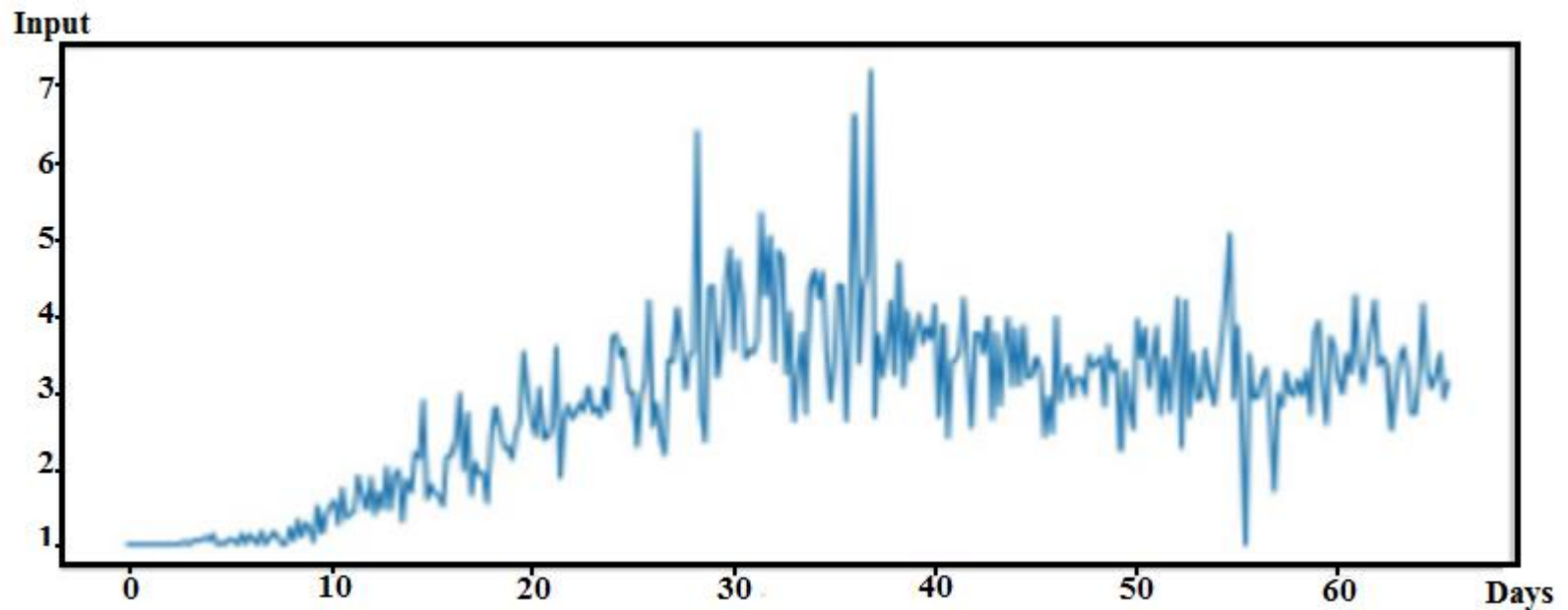
Les données temporelles représentent des valeurs prises par des variables mesurées à des intervalles de temps équidistants (jour, mois, trimestre, année) et ordonnés dans le temps, et cet ordonnancement a une signification que l'on ne peut ignorer.

On représente habituellement une série temporelle X_t ($1 \leq t \leq n$) à l'aide d'un graphique avec en abscisse les dates et en ordonnée les valeurs observées.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

Exemple de série temporelle:



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

La météorologie, la géophysique et l'astrophysique, collectent des observations qui peuvent être représentées sous forme de séries chronologiques.

La température d'une région donnée, la consommation d'électricité sont des exemples de séries temporelles ou de séries chronologiques.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

Données transversales: Données relatives à plusieurs variables à un temps donné.

Il existe un potentiel de corrélation entre les observations à des périodes adjacentes de données des séries chronologiques.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

Modèles de séries temporelles

Les modèles de séries chronologiques:

- Univariées
- Multivariées.

Les modèles de séries chronologiques **univariées** sont des modèles utilisés lorsque la variable dépendante est une série chronologique unique $X = \{x_1; x_2; \dots ; x_m\}$, m : nombre d'observations.

Exemple: Fréquence cardiaque d'un individu par minute.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

Modèles de séries temporelles

Des modèles de séries temporelles **multivariées** : Cas de plusieurs variables souvent dépendantes.

En plus de dépendre de leurs propres valeurs passées, chaque série peut dépendre des valeurs passées et présentes des autres séries.

Une série temporelle **multivariée** avec p variables (X_1, \dots, X_p) et n périodes de temps $(t = 1..n)$.

	X_1	\dots	X_p
1	x_{11}	\dots	x_{p1}
\vdots	\vdots	\ddots	\vdots
n	x_{1n}	\dots	x_{pn}

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Données temporelles

Domaines d'application

- Analyse financière : Elle comprend la prévision des ventes, l'analyse des stocks, l'analyse du marché boursier, l'estimation des prix.
- Analyse météorologique : Elle comprend l'estimation de la température, le changement climatique, la reconnaissance des changements saisonniers, les prévisions météorologiques.
- Analyse des données réseau : Cela comprend la prédiction de l'utilisation du réseau, la détection d'anomalies ou d'intrusions, la maintenance prédictive.
- Analyse médicale et biologie : suivie des évolutions des pathologies, analyse d'électroencéphalogrammes et d'électrocardiogrammes, suivi de maladies cancéreuses.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

- Mesure de similarité en fonction des invariances nécessaires par rapport à l'application :
 - facteur d'échelle en amplitude ou en temps,
 - décalage d'amplitude,
 - déphasage temporel global, ou distorsions temporelles locales,
 - occultations de mesures.
- Algorithme de regroupement de séries chronologiques.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Les principales mesures de similarité

Distances de Minkowski (p-normé)

L'approche la plus simple pour définir la similitude entre deux séquences est de ranger chaque séquence dans un vecteur et puis d'employer une p-distance pour définir la mesure de similarité.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Les principales mesures de similarité

Distances de Minkowski (p-normé)

Soient Q et C deux séries: $Q = q_1, q_2, \dots, q_m$ et $C = c_1, c_2, \dots, c_m$

La distance de Minkowski est définie par:

$$Distance_p(Q, C) = \sum_{i=1}^m ((q_i - c_i)^p)^{1/p}$$

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Les principales mesures de similarité

La similarité (Similitude (Q, C)) de Minkowski est égale à:

$$Similitude(Q, C) = \frac{1}{Distance_p(Q, C)}$$

Si $p=1$, on utilise la distance de Manhattan:

$$Distance_{Manhattan}(Q, C) = \sum_{i=1}^m |q_i - c_i|$$

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Les principales mesures de similarité

La similarité (Similitude (Q, C)) de Minkowski est égale à:

$$Similitude(Q, C) = \frac{1}{Distance_p(Q, C)}$$

Si $p=1$, on utilise la distance Euclidienne:

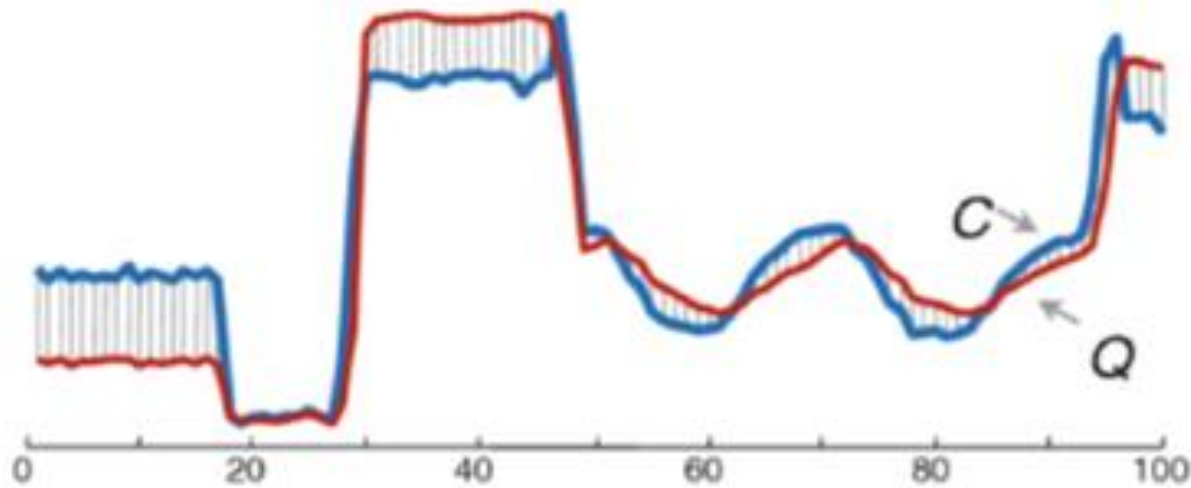
$$Distance_{Euclidienne}(Q, C) = \sqrt{\sum_{i=1}^m (q_i - c_i)^2}$$

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Exemple de calcul de la distance Euclidienne entre deux séries Q et C:



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.

Les points similaires sur les deux séries sont écartés par le temps.

Ainsi, les distances euclidiennes s'élargissent et par conséquent, les séries sont jugées non similaires.

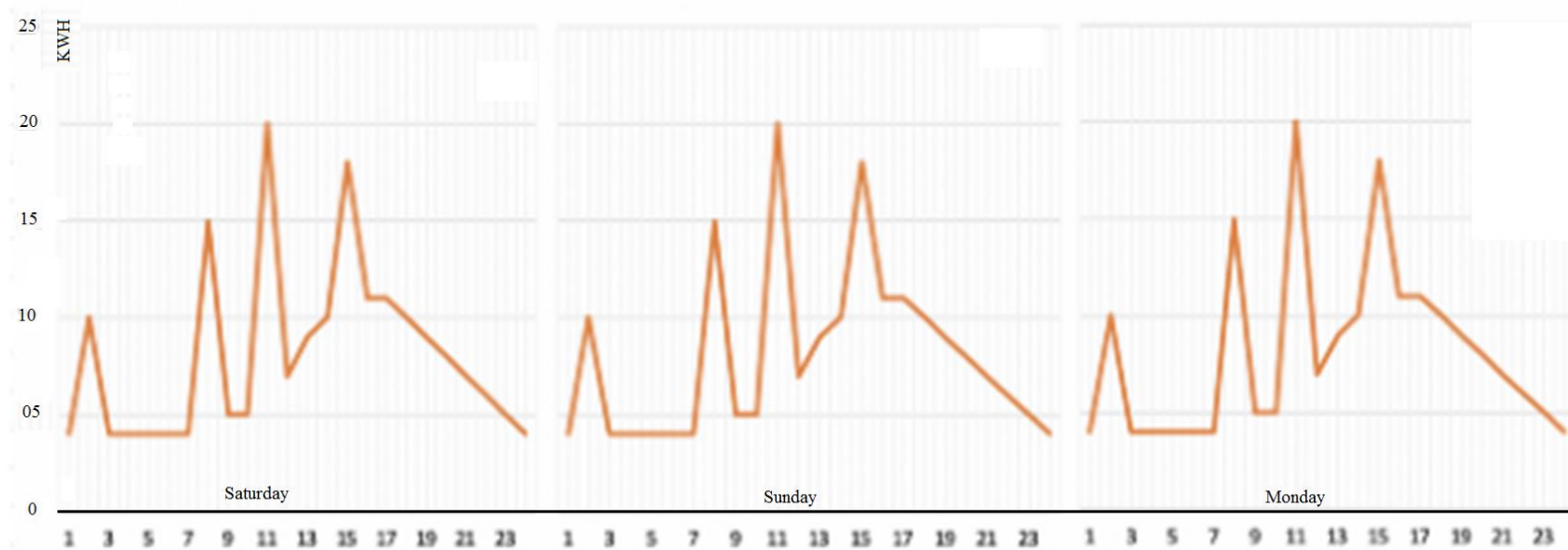
On utilisera ainsi la mesure DTW (Dynamic Time Warping)

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.

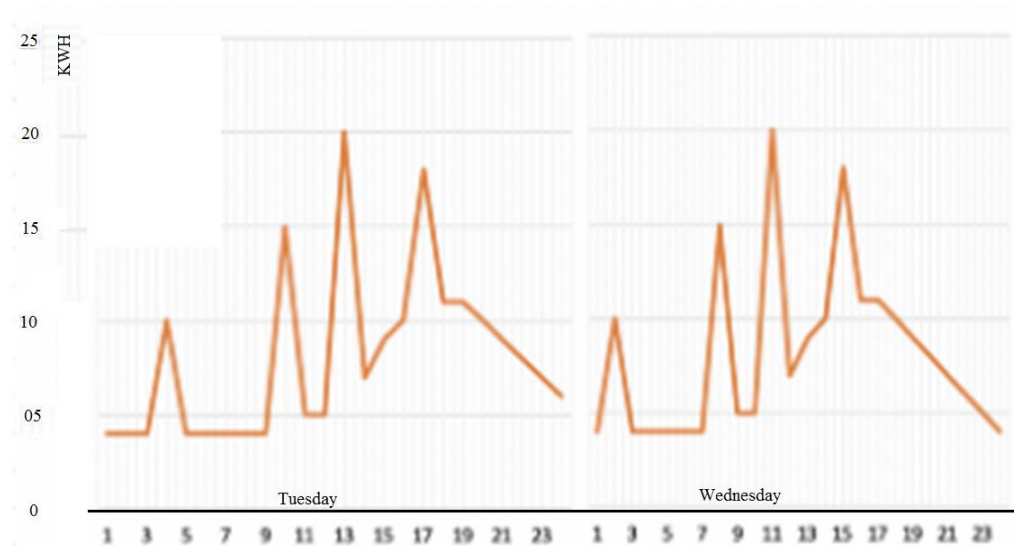


Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.

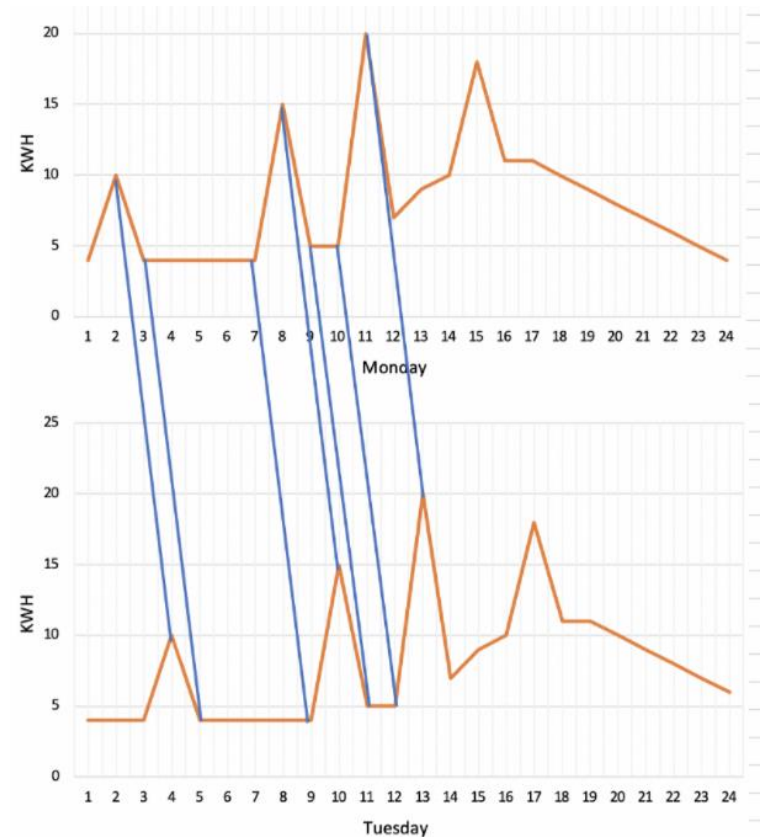


Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

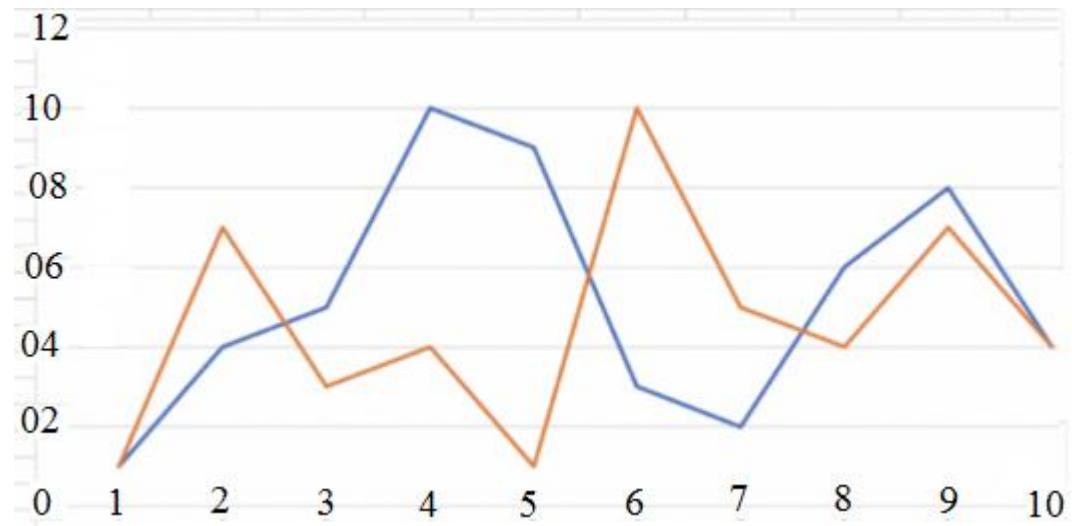
5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.

Série (P) : 1,4,5,10,9,3,2,6,8,4

Série 2 (Q): 1,7,3,4,1,10,5,4,7,4



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

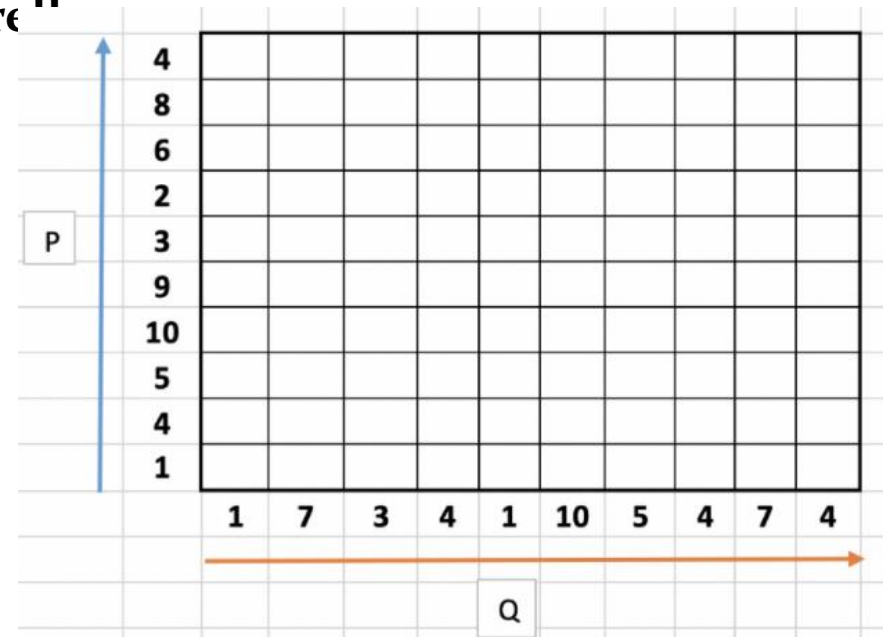
Clustering classique de séries temporelles

Cas de séries non synchronisées.

Série (P) : 1,4,5,10,9,3,2,6,8,4

Série 2 (Q): 1,7,3,4,1,10,5,4,7,4

Création de la matrice M de coûts



$$M(i, j) = |P(i) - Q(j)| + \min (M(i-1, j-1), M(i, j-1), M(i-1, j))$$

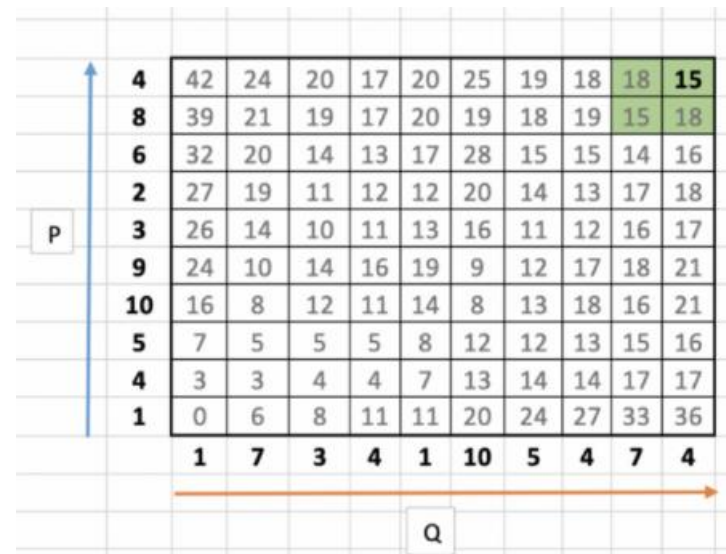
Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.

Le chemin de traversée (traversal path) est identifié sur la base du voisin avec une valeur minimale.



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

Cas de séries non synchronisées.

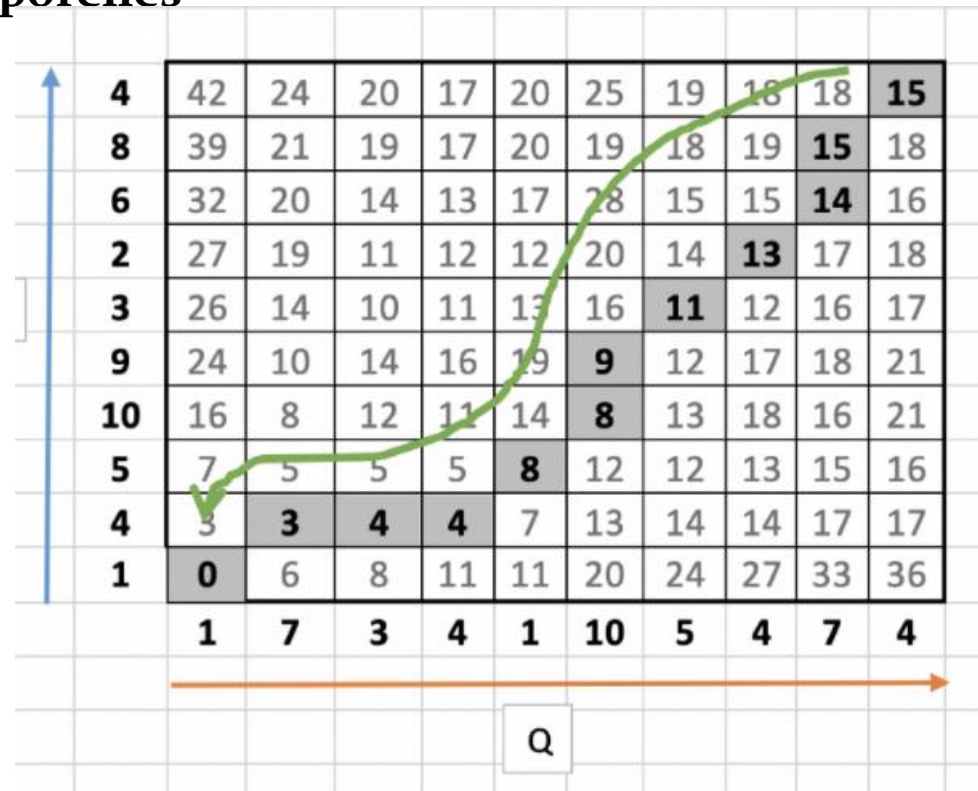
Notons ce chemin par d .

$d = [15, 15, 14, 13, 11, 9, 8, 8, 4, 4, 3, 0]$

Final Distance Calculation

Time normalised distance , D

$$D = \frac{\sum_{i=1}^k d(i)}{\sum_{i=1}^k k} = 8.63$$

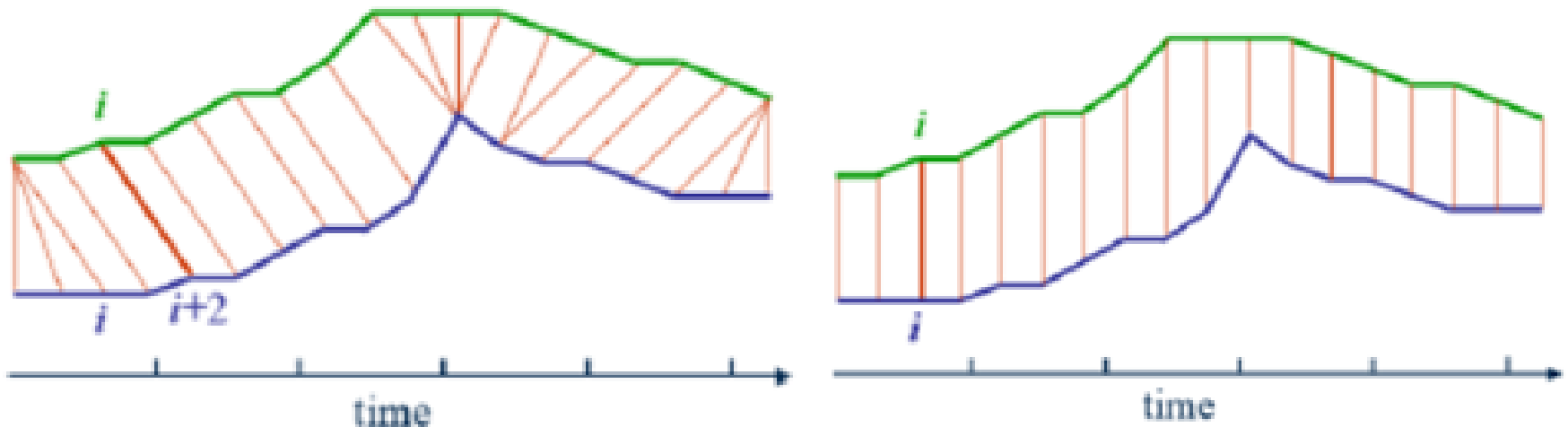


Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.2 Clustering de données temporelles

Clustering classique de séries temporelles

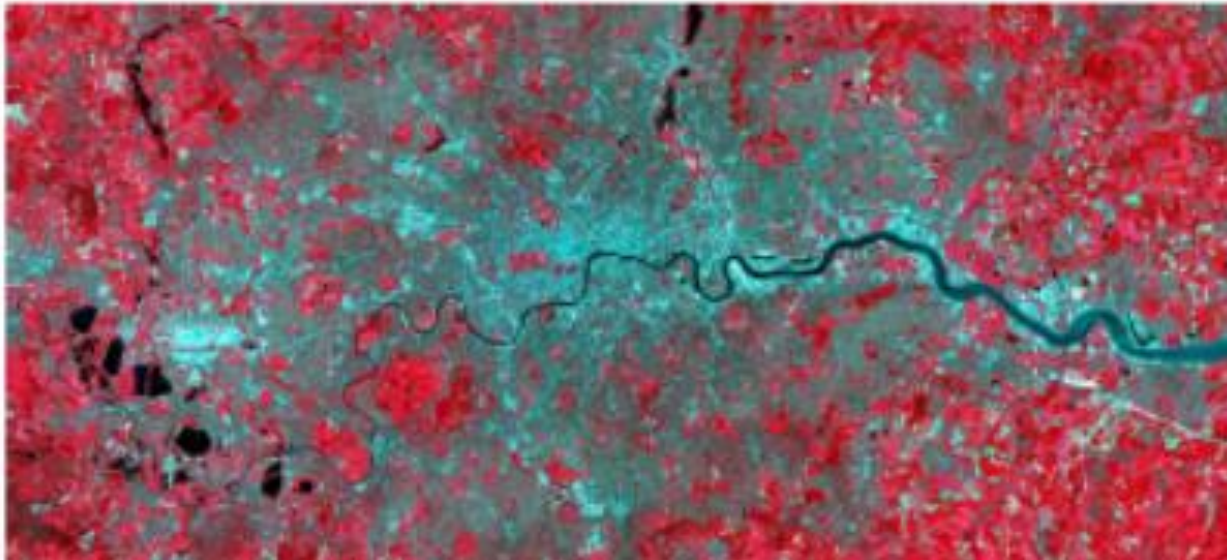
Cas de séries non synchronisées.



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Une image de télédétection est une image d'une scène captée par un système aérien ou satellitaire.



Exemple d'image de télédétection

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Résolution d'image satellitaire :

- spatiale,
- spectrale,
- radio-métrique
- temporelle.

Résolution spatiale : taille de la surface en mètres couverte par un pixel.
la longueur de côté de la surface varie de moins de 1 à 1 000 mètres.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Résolution spectrale : liée au nombre de bandes spectrales (du bleu à l'infrarouge).

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Résolution radiométrique : relative au nombre d'intensités différentes de luminosité que le capteur est capable de distinguer.

Une plage dynamique plus large pour un capteur permet de discerner plus de détails dans l'image.

Le capteur Landsat 7 enregistre des images 8 bits ; il peut ainsi mesurer 256 valeurs de gris uniques tandis qu'Ikonos-2 a une résolution radiométrique de 11 bits (2048 valeurs de gris).

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Résolution temporelle : réfère au temps écoulé entre des images consécutives du même

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Clustering

Le clustering: Regrouper les pixels similaires en régions homogènes afin de réaliser une cartographie du sol ou de la scène télé détectée.

Approche basée pixels: utilisant uniquement les informations spectrales ou radiométrique disponibles pour le pixel individuel.

L'algorithme de classification va générer une classe pour chaque pixel individuel d'une image pour chaque bande spectrale.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Clustering

A noter que:

- Les objets même élémentaires sont plus gros que les pixels
- Pour ces objets, la radiométrie ne suffit plus, les informations des pixels environnants, peuvent aider à identifier correctement la classe du pixel cible.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

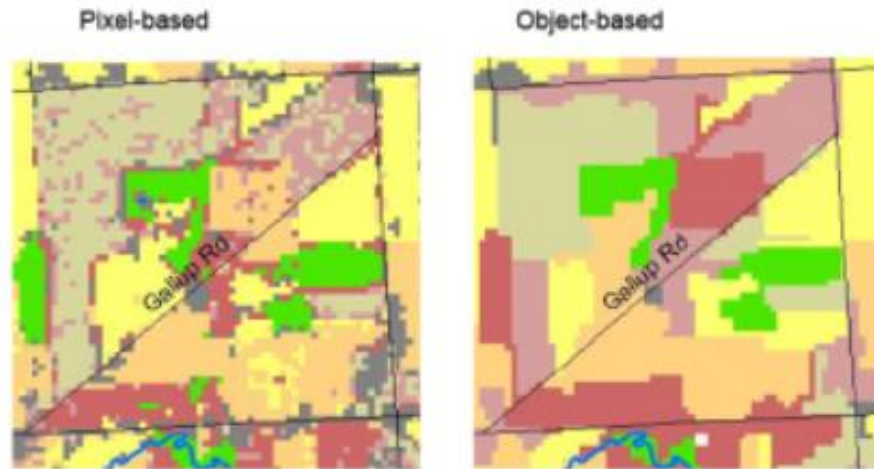
5.3 Cas d'images de télédétection

Clustering

Approche basée objets/régions

Le terme "objet" désigne le groupe de pixels.

L'objectif d'une classification basée sur les objets est de découper (segmenter) l'image en segments internes homogènes de taille variable.



Résultats des deux approches sur une image

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Clustering

Clustering d'une série temporelle d'images de télédétection.

Le clustering sur une série temporelle d'images de télédétection à pour but de regrouper les pixels évoluant dans le temps de manière similaire.

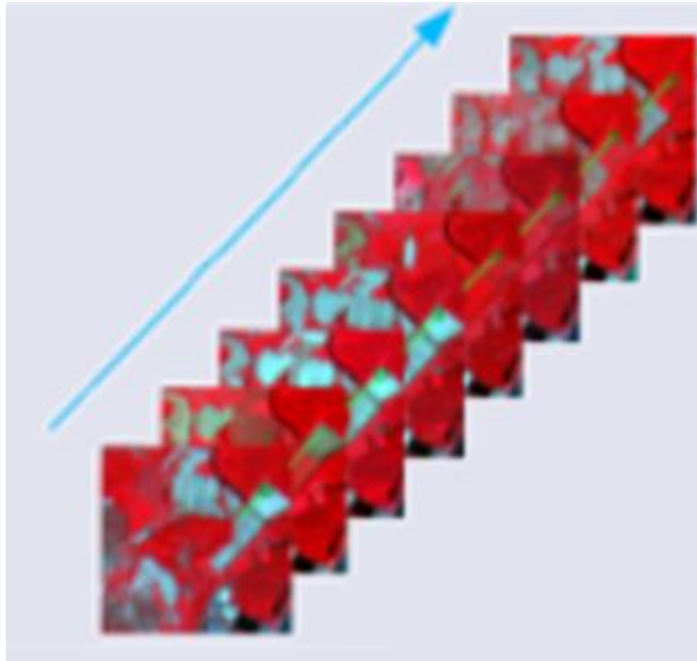
La finalité est de pouvoir suivre l'évolution dans le temps d'une surface terrestre.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Clustering

Clustering d'une série temporelle d'images de télédétection



Série d'images segmentées

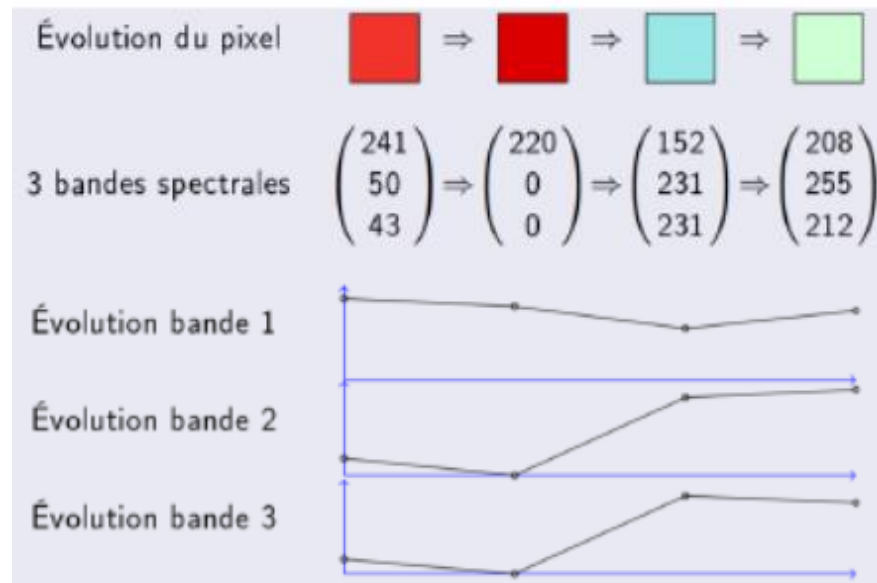
Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Clustering

Clustering d'une série temporelle d'images de télédétection

On construit une séquence par pixel. Une séquence est formée de la suite des états pris par un pixel au cours du temps.



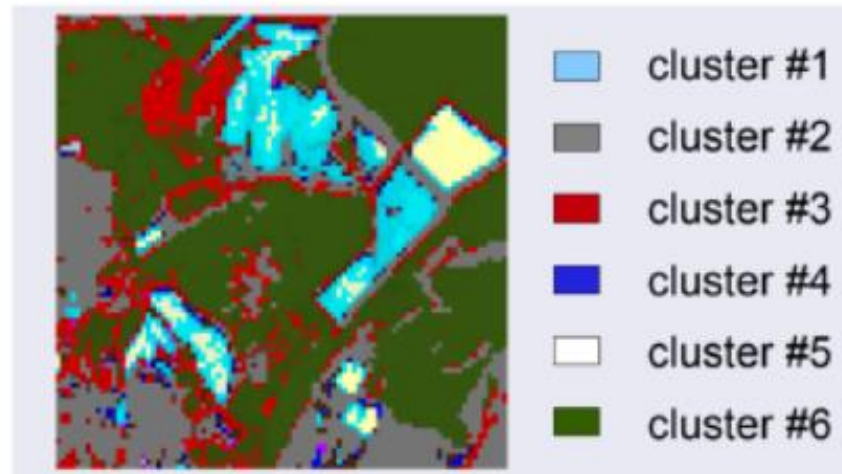
Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.3 Cas d'images de télédétection

Clustering

Clustering d'une série temporelle d'images de télédétection

En utilisant une mesure de similarité entre les séquences des pixels, nous appliquons un algorithme de clustering sur ces séquences.



Resultat d'un clustering basé pixel sur une série d'images de télédétection

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.4 Mise à l'échelle multidimensionnelle (MDS: Multi Dimensional Scaling)

Position du problème:

Construction d'un système de représentation des individus à partir d'une matrice de distances (similarité).

Rendre compte de leurs positions relatives dans un repère euclidien.

MDS: Cette méthode est fondée sur les travaux de Shepard et Kruskal.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Saporta, 2006, page 428 ; restreint aux véhicules étrangers

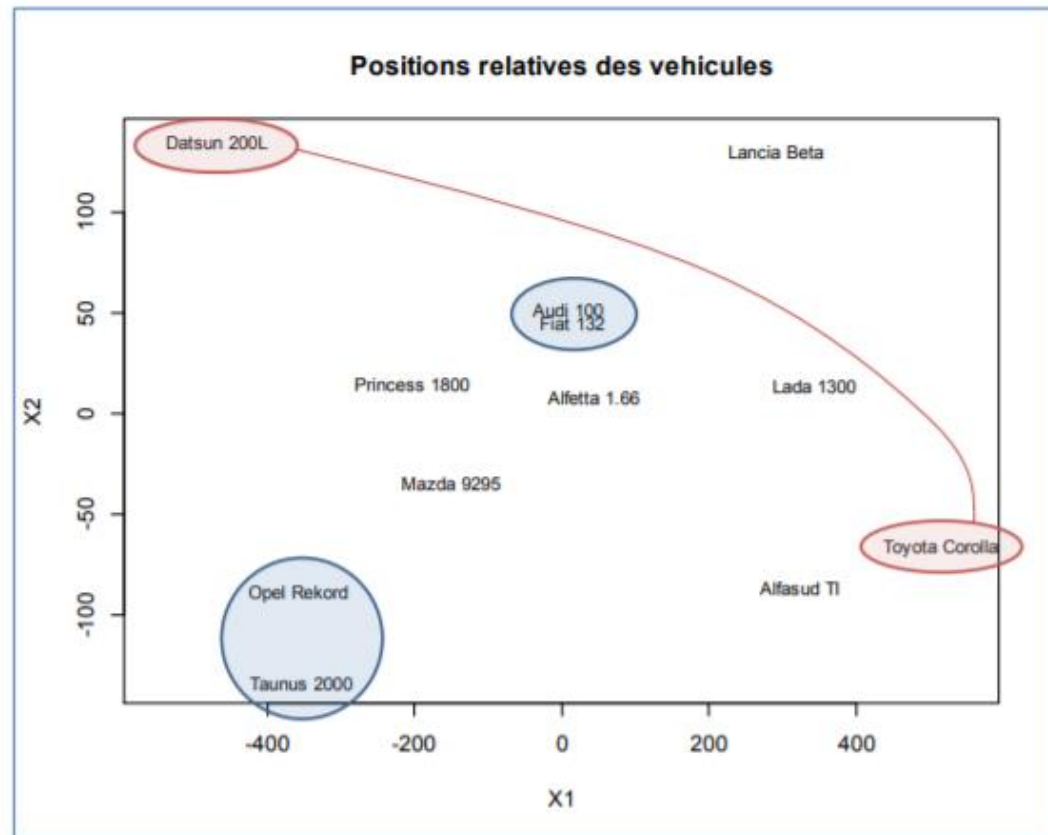
Position du problème: Lecture : bleu, proches ; rouge : éloignés

	Toyota Corolla	Lada 1300	Alfasud TI	Lancia Beta	Mazda 9295	Fiat 132	Alfetta 1.66	Princess 1800	Audi 100	Taurus 2000	Opel Rekord	Datsun 200L
Toyota Corolla												
Lada 1300	190.2											
Alfasud TI	195.3	106.0										
Lancia Beta	299.3	130.1	219.8									
Mazda 9295	667.2	497.4	477.9	472.4								
Fiat 132	513.6	331.6	336.1	289.8	184.9							
Alfetta 1.66	477.8	301.1	294.5	275.6	204.3	51.0						
Princess 1800	722.4	546.2	536.4	507.6	72.0	220.9	251.8					
Audi 100	521.1	339.9	346.7	295.3	184.2	35.2	73.8	217.2				
Taurus 2000	870.9	712.1	678.2	696.3	225.1	408.8	423.9	211.7	407.5			
Opel Rekord	872.3	708.3	680.5	684.2	213.3	395.0	414.8	187.4	391.8	47.8		
Datsun 200L	1004.6	821.4	822.8	760.4	360.5	492.0	530.4	292.9	486.5	292.3	251.7	

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Position du problème:



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Position du problème:

Soit $\delta_{ii'}$ la distance entre deux objets i, i' .

Les positions de ces objets sur une espace 2D sont $p_i(x_i, y_i)$ et $p_{i'}(x_{i'}, y_{i'})$

Il est souhaitable d'avoir

$\delta_{ii'}$ est égale $= \hat{\delta}_{ii'}$ où $\hat{\delta}_{ii'}$ est la distance entre p_i et $p_{i'}$

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Objectifs:

- Se rendre compte visuellement des proximités
- Identifier, interpréter les dimensions qui permettent de discerner les objets
- Se placer sur un espace de représentation qui permet d'utiliser les techniques d'apprentissage automatique.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

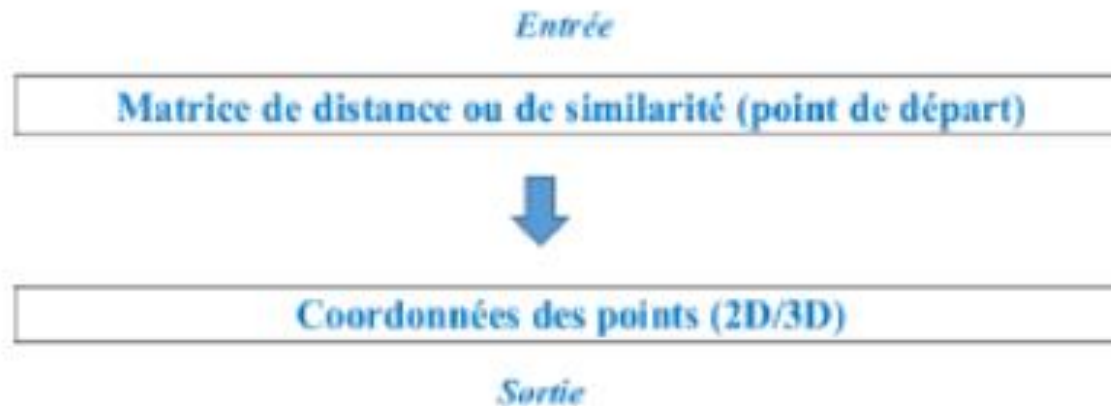
L'échelle multidimensionnelle est une représentation visuelle des distances ou des différences entre les ensembles d'objets.

Étant donné une matrice carré de similitudes perçues entre un ensemble d'objets, MDS trace les objets sur une carte de telle sorte que:

- les objets qui sont perçus comme très similaires soient placés les uns à côté des autres sur le graphique
- les objets qui sont perçus comme très différents les uns des autres sont placés loin les uns des autres sur la carte.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

MDS est aussi utilisé pour servir de technique de réduction de dimension pour des données de grande dimension.

Les données réduites en dimension conservent des propriétés similaires. Par exemple, deux points de données qui sont proches l'un de l'autre dans un espace de grande dimension seront également proches l'un de l'autre dans un espace de faible dimension.

Le multidimensionnel est due au fait que nous sommes pas limités à des graphiques ou des données à deux dimensions. Des tracés en trois dimensions, en quatre dimensions et plus sont possibles.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Mathématiquement, si nous traitons une mesure comme une distance, nous armons que nos mesures ont des propriétés métriques:

- Identité : La distance entre un l'objet et lui-même est de 0.
Alternativement, deux choses qui ont une mesure de distance de 0 sont identiques.
- Symétrie : La distance entre A et B est la même que la distance entre B et A
- Inégalité de triangle : La distance entre A et C doit être inférieure ou égale à la distance entre A et B plus la distance entre B et C.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Exemple du MDS classique : distances entre 10 grandes villes américaines.

	Atl	Chi	Den	Hou	LA	Mia	NY	SF	Sea	DC
Atl	0	587	1212	701	1936	604	748	2139	2182	543
Chi	587	0	920	940	1745	1188	713	1858	1737	597
Den	1212	920	0	879	831	1726	1631	949	1021	1494
Hou	701	940	879	0	1374	968	1420	1645	1891	1220
LA	1936	1745	831	1374	0	2339	2451	347	959	2300
Mia	604	1188	1726	968	2339	0	1092	2594	2734	923
NY	748	713	1631	1420	2451	1092	0	2571	2408	205
SF	2139	1858	949	1645	347	2594	2571	0	678	2442
Sea	2182	1737	1021	1891	959	2734	2408	678	0	2329
DC	543	597	1494	1220	2300	923	205	2442	2329	0

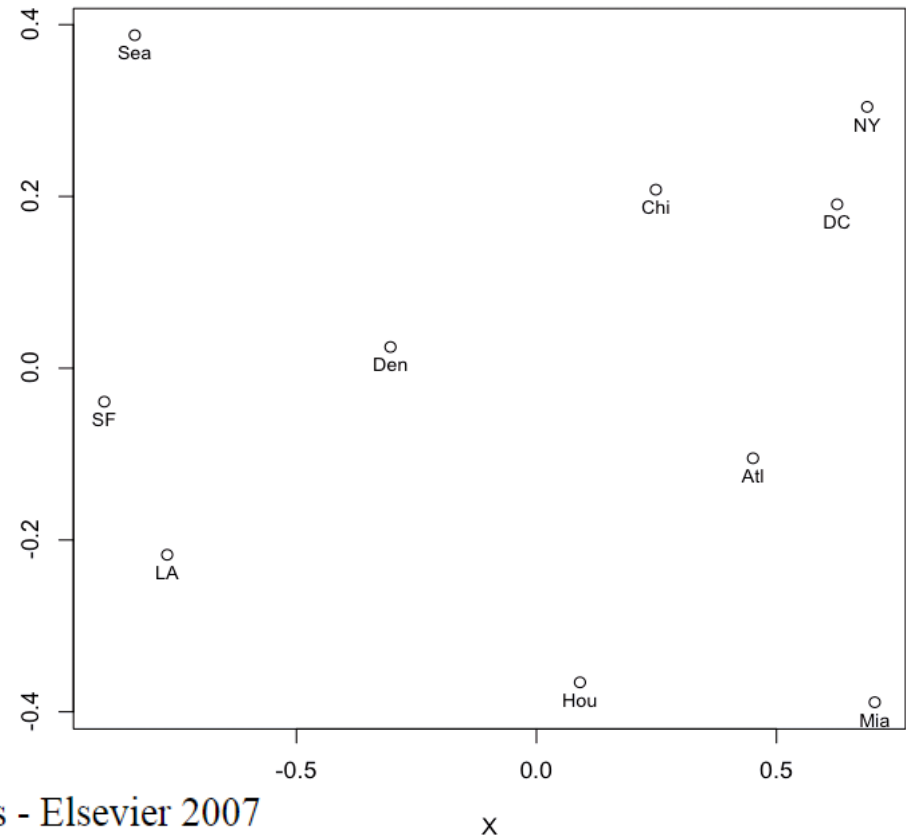
Source : Handbook of Statistics - Elsevier 2007

Exemple : matrice de distance entre les villes

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

MDS produit cette carte :



Source : Handbook of Statistics - Elsevier 2007

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Stress

Le stress est une mesure de qualité de l'ajustement que MDS tente de minimiser, elle varie entre 0 et 1, et est basée sur les différences entre les distances prévues et réelles. La forme générale de ces fonctions est la suivante :

Dans son article MDS, Kruskal a écrit que les ajustements proches de zéro sont excellents, alors que tout ce qui dépasse 0,2 devrait être considéré comme médiocre.

STRESS > 0.20 : mauvais

0.10 < STRESS < 0.20 : passable

0.05 < STRESS < 0.025 : bien

STRESS < 0.025 : excellent

STRESS = 0 : parfait.

Kruskal, J. B. (1964). "Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis". *Psychometrika*. 29 (1): 1–27.
doi:10.1007/BF02289565

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Stress

STRESS > 0.20 : mauvais

0.10 < STRESS < 0.20 : passable

0.05 < STRESS < 0.025 : bien

STRESS < 0.025 : excellent

STRESS = 0 : parfait.

Des auteurs plus récents suggèrent d'évaluer le stress en fonction de la qualité de la matrice de distance et du nombre d'objets dans cette matrice.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Diagramme de Shepard

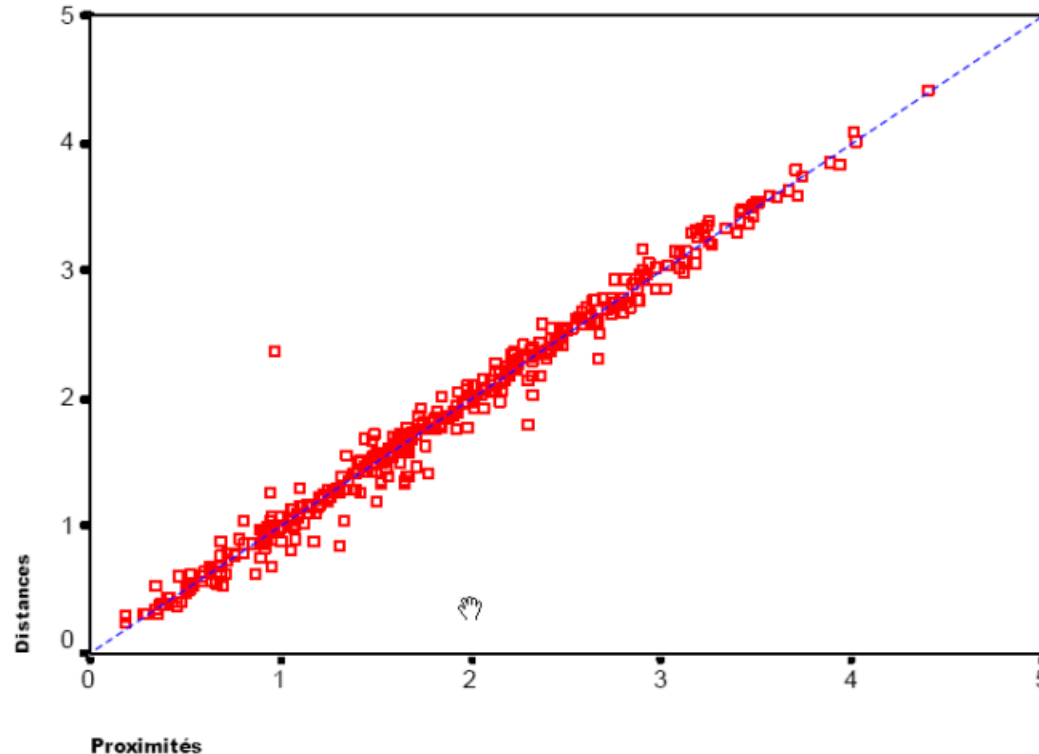
Un diagramme de Shepard compare la distance entre les points de données avant et après leur transformation sous forme de nuage de points.

Les diagrammes Shepard peuvent être utilisés pour évaluer la qualité d'ajustement des techniques de réduction des données telles que l'analyse des composants principaux (ACP), la mise à l'échelle multidimensionnelle (MDS).

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

Diagramme de Shepard



Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

MDS: Les algorithmes:

MDS est une famille d'algorithmes différents, chacun conçu pour arriver à une configuration optimale de faible dimension ($p = 2$ ou 3).

MDS classique : Le MDS classique est l'approche simple et rapide, il génère une carte en réduisant l'erreur entre les distances réelles entre les éléments et la distance cartésienne entre les éléments sur la carte.

Elle tente de reproduire la métrique ou les distances d'origine.

Chapitre 5: VISUALISATION DE DONNEES HAUTE DIMENSION

5.1 Mise à l'échelle multidimensionnelle (MDS)

MDS non métrique: Dans le MDS non métrique, on cherche à préserver l'ordre des proximités et non leurs valeurs absolues ou relatives. Autrement dit, le but est de représenter les distances entre les objets, en respectant l'ordre entre les proximités plutôt que leurs valeurs exactes.

MDS isométrique: La carte résultant de la MDS classique est prise, puis ajustée de sorte que les distances sur la carte entre des paires d'éléments apparaissent dans le même ordre décroissant que les données d'origine. Exemple: différence de la valeur nutritive entre les aliments, où un classement des distances est plus important qu'une coordonnées d'unité spécifique.