

State of the art

# **How to evaluate the faithfulness of visual data projection ?**

June 16, 2022

### Students :

|        |          |                                      |
|--------|----------|--------------------------------------|
| Cyril  | GARDENAT | cyril.gardenat9@etu.univ-lorraine.fr |
| Nathan | METZGER  | nathan.metzger6@etu.univ-lorraine.fr |

### Tutors :

Lydia BOUDJELOUD-ASSALA

**Keywords:** projection; méthodes; critères; évaluation; qualité; réduction de dimension

### Abstract :

La visualisation de données est une façon de représenter des jeux de données par le biais de graphiques. C'est une façon très efficace de communiquer sur des données, d'autant plus lorsque que celles-ci sont complexes et présentes en grande quantité. Dans ce rapport, nous nous sommes concentrés sur une approche de la visualisation de données : la projection. Ainsi, différentes méthodes de projection (linéaires et non linéaires, supervisées et non supervisées) ont été présentées de façon non exhaustive. Une perte de donnée étant inévitable lors de la projection de grands jeux de données, la problématique de la qualité de cette dernière se pose. De ce fait, nous avons également cité les biais qui peuvent découler d'une projection ainsi que les critères qui permettent de juger de sa qualité. Enfin nous avons présenté le programme que nous avons proposé pour évaluer la qualité des différentes méthodes de projection.

# Contents

|   |    |
|---|----|
| Introduction . . . . .                                    | 1  |
| 1 Les méthodes de projection . . . . .                    | 4  |
| 1.1 Analyse en Composantes Principales(ACP) . . . . .     | 5  |
| 1.2 Multidimensional Scaling (MDS) . . . . .              | 5  |
| 1.3 ClassNerv . . . . .                                   | 6  |
| 1.4 t-SNE . . . . .                                       | 6  |
| 1.5 Autres méthodes de projection non linéaires . . . . . | 8  |
| 1.6 Conclusion de la partie . . . . .                     | 8  |
| 2 Les critères d'évaluation des projections . . . . .     | 9  |
| 2.1 Visualisation du stress . . . . .                     | 9  |
| 2.2 Les critères de mesure . . . . .                      | 9  |
| 2.3 Conclusion de la partie . . . . .                     | 12 |
| Réalisation . . . . .                                     | 14 |
| 3 Implémentation . . . . .                                | 14 |
| 3.1 Méthodologie . . . . .                                | 14 |
| 3.2 Jeux de données utilisés . . . . .                    | 14 |
| 3.3 Fonctionnement du programme . . . . .                 | 15 |
| 3.4 Expérimentation . . . . .                             | 16 |
| 4 Analyse des résultats . . . . .                         | 16 |
| 4.1 Analyse de chaque critères . . . . .                  | 16 |
| 4.2 Discussion des résultats . . . . .                    | 19 |
| 5 Conclusion . . . . .                                    | 19 |

## Introduction

La représentation visuelle des grands jeux de données est un enjeu de plus en plus important dans de nombreux domaines scientifiques. En effet, une immense partie des analyses de données se fait dans des domaines de recherches dont les acteurs ne sont pas des experts en analyse de données. Cela rend donc l'analyse compliquée et peut engendrer certains biais lors de leur interprétation [1]. C'est pour cela, que le domaine de la visualisation des données se doit d'être accessible à ce genre de profil.

La visualisation permet aux utilisateurs d'avoir un aperçu global des données et de choisir en conséquence quel type de traitement et quel paramétrage utiliser dans leurs analyses. Elle a donc pour rôle d'être une interface interactive pour représenter et naviguer à travers les données extraites, et ce dans le but de les rendre compréhensibles et par conséquent exploitables[2]. C'est exactement ce vers quoi tend ce domaine : “ *L'analyse visuelle des données combine techniques automatiques d'analyse et visualisation interactive dans le but de pouvoir comprendre, réfléchir et prendre des décisions de la bonne façon en se basant sur des jeux de données grands et complexes.* ” [3]. Différentes techniques de visualisation existent [4], mais c'est sur la projection que nous allons nous concentrer.

Dans une projection, toutes les données sont projetées selon une projection linéaire ou non des axes dans un plan en deux dimensions. Cela permet de représenter les données sous la forme d'un nuage de points on parle de représentation dispersée (ou *scattered representation*) en essayant de respecter au mieux la structure des données [4].

Dans une projection, les axes ne jouent pas forcément un rôle: c'est principalement la distance entre les points qui donne un sens et une interprétabilité à la projection[5]. Le principe de base d'une lecture projection est le suivant : les points qui sont proches les uns des autres sont censés être relativement similaires tandis que ceux qui sont éloignés sont censés être différents.

Les projections permettent également d'inférer des propriétés sur les données et de les vérifier. Par exemple, il est possible de vérifier si les clusters sont séparables linéairement (notamment à l'aide d'une projection linéaire) [4] .

Acquérir les données et les traiter sont deux étapes qui sont de plus en plus simples. Le défi de la projection est de représenter fidèlement les jeux de données. En effet, dans une projection les jeux de données peuvent atteindre jusqu'à des dizaines de milliers de dimensions et doivent être réduits en deux ou trois dimensions pour être visualisables et compréhensibles par l'homme. Néanmoins, représenter fidèlement une distance euclidienne après réduction et projection est tout bonnement impossible car ce processus de réduction de dimension a un coût une : perte d'informations. Celle-ci est causée par l'étirement ou la compression de la plupart des distances entre paires, car la dimension de l'espace couvert par les données d'origine est généralement supérieure à la dimension de l'espace de projection. Cela cause l'apparition d'artefacts dans la projection. Ces artefacts compromettent ainsi grandement la fiabilité et la bonne interprétation de ces graphiques[6].

Par conséquent, la mesure et la visualisation de ces distorsions appelées artefacts sont cruciales pour l'analyse. Elles se doivent d'être analysées afin de détecter quelles distances entre paires ont été préservées, et donc d'évaluer si les caractéristiques observées dans l'espace de projection sont des images fidèles des caractéristiques de l'espace originel, ou simplement des artefacts de la projection.

Les artefacts sont définis comme des points mal placés(en comparaison aux points d'origine). Ils sont dus au processus de réduction qui peine à respecter fidèlement les distances. Il existe deux types d'artefacts : les artefacts géométriques et les artefacts topologiques [1].

- Les artefacts géométriques sont causés par de légères distorsions des distances, dans ce cas précis les distances sont fausses, mais le voisinage des points est conservé, l'interprétation n'est donc pas (ou très peu) perturbée.

- Les artefacts topologiques sont causés par des distorsions trop importantes des distances. Par conséquent, cela engendre forcément des problèmes d'interprétation.

Parmi ces types d'artefacts, nous pouvons distinguer 4 sous-types [7] :

- Les compressions : Les distances par paire de projections sont inférieures aux distances par paire d'origine correspondante, mais la topologie du voisinage d'origine est préservée.
- Les étirements : Les distances par paires projetées sont plus grandes que les distances par paires d'origine correspondantes, mais la topologie du voisinage d'origine est préservée.
- Les collages : Une compression particulière dans laquelle les points très éloignés dans l'espace d'origine deviennent des voisins proches dans l'espace de projection, changeant ainsi radicalement la topologie du voisinage.
- Les déchirements : Un étirement particulier où des points proches dans l'espace d'origine sont projetés loin les uns des autres, changeant drastiquement la topologie du voisinage.

Il y aurait deux causes principales à l'origine de ces artefacts [7] :

- Les causes structurelles(ou intrinsèques) : Les structures géométriques et topologiques des collecteurs d'origine et celles de l'espace de projection ne sont pas compatibles, de sorte que la projection ne peut se faire sans modifier les distances par paires entre les données.
- Causes techniques(ou extrinsèques) : Si la technique de projection est non linéaire, elle peut créer des artefacts qui n'existent pas à l'optimum global et donc modifier la forme initiale.

Les artefacts ayant des causes techniques peuvent être annulés alors que ceux ayant des causes structurelles ne le peuvent pas. Cependant, dans la pratique, il n'est guère possible de distinguer les deux causes d'artefacts dans le cas des projections non-linéaires.

Pour tenter de pallier les difficultés induites par les artefacts, des techniques ont été mises en place pour assister l'utilisateur dans la visualisation. Parmi ces techniques, nous pouvons noter l'utilisation d'échelles colorimétriques[8], qui permettent de mesurer les différents artefacts, mais ne permettent pas de mettre en valeur les clusters cachés ni de donner une idée claire ou de préciser de la qualité de la projection. Dans ce cas, il est compliqué d'exploiter une projection de  $n$  dimensions de façon fiable sans prendre en compte les artefacts qu'elle présente[7]. Dans l'article "Visualizing Dimensionality Reduction Artifacts : An evaluation"(Heulot) une technique colorimétrique est utilisée pour voir s'il est possible de surmonter les difficultés dues aux artefacts dans d'interprétation d'échelles multidimensionnelles. Ils ont donc mis en place une méthode s'appelant "ProxiViz". Grâce à celle-ci, sur un jeu de test, les informations locales d'un item sont beaucoup mieux représentées. Cette méthode permettrait de détecter plus facilement les outliers et les clusters. Elle pourrait donc, en partie, aider les non spécialistes en analyse de données à mieux distinguer certaines informations. Même si encore une fois cela est limité par la qualité de l'écran, le contraste et la vision de l'humain qui peut peiner à distinguer les nuances entre les couleurs[1] [9] [10]. L'idée a notamment été reprise par Heulot[4], il se base sur le concept de visualisation interactive des proximités afin de mettre en évidence les artefacts. Cela permet de visualiser sur la projection les proximités d'origines des différents points en fonction d'une référence choisie par l'utilisateur.

Comme vous avez pu le comprendre, la qualité d'une projection est primordiale pour une bonne interprétation des données. La qualité de la projection dépend du jeu de données et de la méthode utilisée pour le réduire et le projeter. Par conséquent, il est parfaitement logique de se demander comment évaluer la qualité d'une projection et comment choisir la méthode de projection en fonction d'un jeu de données. Pour tenter de répondre à cette question nous avons mis en place un programme permettant d'évaluer les différentes méthodes de projection pour un même jeu de

données. Dans ce rapport nous parlerons dans un premier temps des différentes méthodes de projections que nous allons étudier. Puis nous expliciterons les critères que nous allons utiliser pour les évaluer. Enfin nous parlerons en détails du programme que nous avons codé permettant d'évaluer ces différents types de projections.

# 1 Les méthodes de projection

Lorsque l'on souhaite représenter en deux ou dimensions de grands jeux de données multidimensionnels (au-delà de vingt dimensions [11][4]), il faut passer par une réduction de dimension. Cette étape intervient dans le processus de projection et est indispensable dans ce cas précis où il y a un nombre important de dimensions (plus d'une vingtaine [4]). Néanmoins, une réduction implique inexorablement une perte d'information et donc une perte de précision, ce qui peut biaiser la projection ainsi que son interprétation. Il est donc important, malgré des pertes inévitables, d'essayer de conserver au maximum les structures et les données utiles. Le facteur numérique est donc un facteur clé, qui se doit d'être compris et maîtrisé.

Dans l'espace des données, il est possible d'avoir recours à différents types de mesures (différents critères) pour déterminer à quel point deux points de la projection sont similaires. Cela dépend grandement de la sémantique sous-jacente à la notion de similarité, celle-ci étant liée au domaine scientifique d'où proviennent les données et à leur nature [4]. Comme mentionné plus haut, la projection se base sur un encodage des similarités par le biais de la variable de position: plus les objets se ressemblent, plus ils seront proches sur la projection, et inversement. Dans une projection, une dissimilarité métrique dans l'espace des données correspond à la définition mathématique d'une distance [12] [4]

$$d : \mathbb{R}^m * \mathbb{R}^m \rightarrow \mathbb{R}$$

, qui pour

$$\forall (u, v, w) \in \mathbb{R}^{3*m}$$

satisfasse les conditions suivantes [4] :

- La définition :  $d(u, v) = 0 \Leftrightarrow u = v$
- La positivité :  $d(u, v) \geq 0$
- La symétrie :  $d(u, v) = d(v, u)$
- L'inégalité triangulaire :  $d(u, v) \leq d(u, w) + d(w, v)$

La distance de Manhattan et la distance euclidienne peuvent subir la domination d'une dimension qui aurait des valeurs réparties sur un spectre plus large que les autres [4]. Pour corriger cela, on normalise ou on pondère les dimensions (il faut cependant faire attention à bien prendre en compte les poids dans l'interprétation de la similarité).

Il est possible de séparer les méthodes de projections en plusieurs catégories. Elles peuvent être non supervisées ou supervisées : c'est-à-dire qu'elles nécessitent l'étiquetage des données et/ou l'intervention de l'utilisateur pour positionner les points sur la projection [13] [4]. Les algorithmes de projections peuvent également être séparés en deux catégories : les algorithmes de projections linéaires et non linéaires. Dans le cas des algorithmes de projection linéaire, l'espace de plus faible dimension (celui de la projection) est obtenu par la combinaison linéaire des dimensions de l'espace de départ.

Les algorithmes de projection non linéaires partent du principe qu'il existe une variété non linéaire de plus faibles dimensions que l'on peut projeter sans trop de déformations dans un espace 2D (ou 3D). Les données sont ainsi projetées en conservant les propriétés de leur variété. Pour cela, elles utilisent des mesures de stress métriques (selon les distances) ou Non-Métriques (selon le rang des distances). Ensuite plusieurs méthodes d'optimisation sont utilisées telles que la descente de gradient, les réseaux de neurones [14] ou bien le placement par force (pour les jeux de données les moins gros) [15]. Cela permet soit de trouver un optimum global (qui préserve la structure globale) ou un optimum local (dans le cas où les méthodes préservent les voisinages locaux). Pour illustrer la

présentation de ces méthodes nous utiliserons le jeu de données GLOBE qui contient 512 données réparties aléatoirement sur la surface d'un globe(donc en trois dimensions). Il est composé de deux classes (bleue et rouges) qui correspondent à deux hémisphères séparés au niveau de l'équateur [16].

### 1.1 Analyse en Composantes Principales(ACP)

Ce type de projection linéaire vise à préserver de la façon la plus optimale la variance dans les données. Cet objectif est atteignable en transformant les variables originelles fortement corrélées en variables non corrélées les unes des autres(également appelées *composantes principales*) qui expliquent la variabilité dans les données. L'ACP mesure l'inertie des données pour réduire au maximum le nombre de variables. Géométriquement parlant, elle a recours à une projection orthogonale des données dans un sous-espace principal obtenu par combinaison des dimensions dans le jeu de données. Le sous-espace est calculé de sorte que le nuage de points obtenu à partir de la projection maximise la variance des données [4].

Cette méthode est beaucoup utilisée dans de nombreux domaines d'application. Elle fournit une pondération des composantes principales, ce qui permet de déterminer les directions de dispersions des données les plus importantes. La projection résultante(Figure 1) est directement interprétable,car les vecteurs propres décrivent la contribution de chaque variable [4]. Néanmoins, elle ne permet pas de séparer des relations non linéaires et est plus susceptible d'introduire du faux voisinage entre les points.

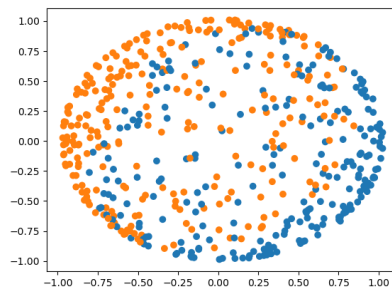


Figure 1: Réduction du dataset Globe par ACP

Cette méthode est très performante en termes de temps de calcul et respecte fortement la distance au voisinage de chaque point, cependant elle ne représente pas correctement les classes. La méthode est implémentée dans la librairie scikit-learn.

### 1.2 Multidimensional Scaling (MDS)

Le principe de fonctionnement de cette méthode est de projeter les données dans un sous-espace linéaire[17]. L'objectif de cette approche est de conserver de la meilleure façon possible les distances au carré provenant de l'espace des données. Pour y parvenir, l'algorithme cherche pour chaque métrique de l'espace des données, une combinaison linéaire optimale. Globalement cette technique a les mêmes défauts que l'ACP quand il s'agit de projeter des données linéairement séparables. Néanmoins, dû au fait que la cette technique prend une matrice de similarité en entrée, les axes de projection de cette technique ne sont pas directement interprétables[4].



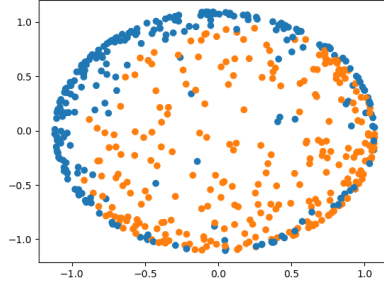


Figure 2: Réduction du dataset Globe par MDS

Cette méthode est performante en termes de temps de calcul et de mémoire. Dans le cas du jeu de données *GLOBE* que nous avons utilisé, les classes sont préservées, néanmoins une d'entre elles est projetée à l'extérieur.

La méthode est implémentée dans la librairie scikit-learn. <https://scikit-learn.org/stable/modules/generated/sklearn.manifold.MDS.html>

### 1.3 ClassNerv

Cette méthode a été présentée lors de la prestigieuse conférence NeurIPS de 2020 [16].

C'est une méthode de projection qui permet de réduire la dimension d'un jeu de données en conservant les classes et leur voisinage. Elle présente les avantages des méthodes supervisées et non supervisées, c'est à dire la combinaison du traitement par classe et par distance au voisinage ce qui permet d'éviter de nombreuses erreurs de classification. On peut y faire varier les paramètres

$$\tau^* = \frac{\tau^\in + \tau^\notin}{2}$$

$$\varepsilon = \frac{\tau^\in - \tau^\notin}{2}$$

$\tau^* \in [0, 1]$  avec 0 pour un minimum de faux voisin et 1 pour un minimum de voisin raté.

$\varepsilon \in [0, 0.5]$  0 pour réduire la supervision, et 0.5 pour l'augmenter.

Par exemple sur le Dataset Globe (deux classes séparées en deux hémisphères, représentation d'origine en 3D), on obtient les résultats suivants :

On peut remarquer que l'on obtient un meilleur résultat avec un  $\varepsilon$  de 0.25 et un  $\tau^*$  de 0.75, c'est-à-dire une supervision moyenne et en privilégiant les faux voisins plutôt que les voisins manqués. Il est également important de noter que cette méthode est très coûteuse en temps de calcul.

Le code Python nécessaire pour l'implémentation de cette méthode est accessible ici : [link](#).

### 1.4 t-SNE

Cette méthode non linéaire se base sur un algorithme de projection: le SNE [18] celui-ci permet de bonnes visualisations, mais a une fonction de coût qui est difficile à optimiser. C'est notamment de ce problème dont le t-SNE s'affranchit en utilisant une fonction coût symétrique et une t-distribution plutôt qu'une distribution Gaussienne. Ces deux changements lui permettent d'enlever les problèmes d'optimisation du SNE. Cette méthode est donc plus simple et plus rapide à exécuter. De plus, les maps produites seraient même de meilleure qualité[19]. Avec le t-SNE, il est possible de restituer fidèlement la structure d'un grand jeu de données au niveau local et au niveau

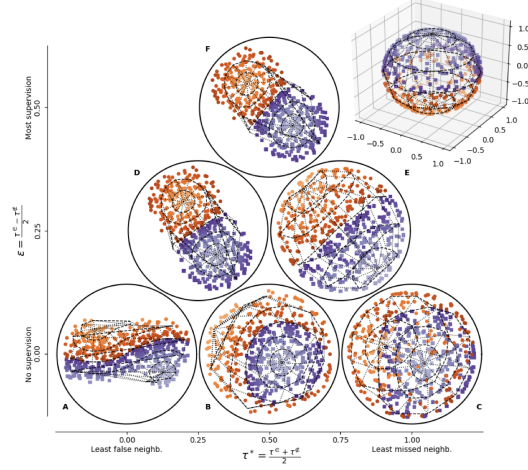


Figure 3: Réduction du dataset Globe par ClassNerV

global. Cette méthode permet en effet de repérer la présence de clusters à différentes échelles. Il est par ailleurs possible de retrouver les clusters “séparés”, et de faire une approximation du spectral clustering avec un certain paramétrage. [20] [21] [22].

La démarche de cette méthode est la suivante : elle commence par calculer la probabilité  $P$  que deux instances puissent être voisines. Ainsi, la valeur de  $P$  sera haute pour des voisins proches et très faible pour des voisins éloignés.

La variance de la distribution Gaussienne est différente à chaque fois, ce qui permet d’enregistrer la densité pour différents voisinages multidimensionnels. Le calcul va s’effectuer de façon itérative jusqu’à ce que la perplexité (au préalable définie par l’utilisateur) soit atteinte.

Une fois, toutes les probabilités calculées le but est de trouver une distribution de probabilité  $Q$  qui représente fidèlement le jeu dans un espace de plus faible dimension. Cette fois-ci, au lieu d’utiliser une distribution Gaussienne, c’est une distribution du  $t$  de Student qui est utilisée, avec un degré de liberté de 1.

Contrairement à  $P$ ,  $Q$  n’est pas paramétré avec une variable de densité de voisinage, par conséquent des voisinages avec des densités très différentes dans l’espace original peuvent être transformés en voisinages de tailles équivalents dans la représentation à plus faible dimension. La recherche de  $Q$  qui représente fidèlement  $P$  se fait en optimisant la fonction coût suivant la divergence de Kull-Leibler des deux distributions. Cette étape est effectuée de façon itérative en faisant une descente de gradient pour un nombre d’itérations déterminé par l’utilisateur [19].

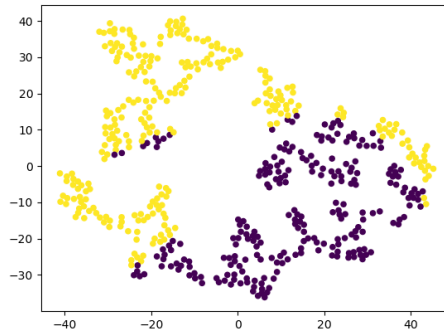


Figure 4: Réduction du dataset Globe par tSNE

Cette méthode est très performante en temps de calcul, mais cela se fait au détriment des résultats qui sont beaucoup moins précis. Les classes sont relativement bien séparées, mais la distance au voisinage n'est absolument pas respectée.

Celle-ci est implémentée dans des projets github par Laurens Van Der Maaten dans une grande diversité de langages (Matlab, Python/Cuda, R, ...).

## 1.5 Autres méthodes de projection non linéaires

La liste des méthodes de projections non linéaires ne se limite évidemment pas à celles citées. Nous pouvons également retrouver le Sammon mapping[23], le *curvilinear components analysis* (CCA)[24], le *Maximum Variance Unfolding* (MVU)[25] [26], ou le *Laplacian Eigenmaps*[27]. Ces techniques sont performantes, mais peinent à être efficaces quand il s'agit de visualiser des jeux de données avec beaucoup de dimensions. Il nous était également compliqué de restituer correctement les structures locales et globales dans une seule projection[28] [19]. C'est en partie pour ces raisons et le manque d'implémentation disponibles qu'elles n'ont pas été retenues pour notre projet.

## 1.6 Conclusion de la partie

Comme nous avons pu le voir, il existe de nombreux algorithmes de projections qui se basent tous sur différentes hypothèses mathématiques, différents critères et différents paramétrages. Ainsi, avec un même jeu de données, il est possible d'obtenir plusieurs projections différentes en fonction des algorithmes choisis (le nuage de points ne sera pas le même). Le choix du type d'algorithme de projection à utiliser dépend de la nature des données que l'on veut étudier, du type de résultat que l'on souhaite observer et de nombreuses caractéristiques propres aux données qui peuvent être très complexes à appréhender sans de bonnes connaissances en projection de données.

Par exemple, nous avons vu pour L'ACP, qu'elle n'est pas stable face aux variations des données. La raison est que l'ajout ou la suppression de quelques instances peut modifier de manière significative les vecteurs et les valeurs propres, affectant ainsi considérablement la cartographie. Même de petites perturbations dans les données peuvent avoir un impact conséquent sur la cartographie finale en raison de l'inversion des vecteurs propres. Ce qui entraîne des configurations très différentes avant et après la perturbation [29]. Le t-SNE quant à lui est basé sur des procédures d'optimisation dont la solution optimale dépend des conditions initiales. En général, des initialisations différentes conduisent à des solutions différentes. La stabilité n'est pas garantie même en fixant la condition initiale, de sorte qu'une modification du nombre d'éléments peut conduire à des agencements sensiblement différents [30]. Quant au MDS, c'est une technique qui repose sur des points de repère ou de contrôle et qui a tendance à être stable tant que les points de contrôle et de repère restent fixes. Comme la position de chaque instance ne dépend que des points de contrôle (points de repère), si ces points ne changent pas, le mappage ne change pas non plus [30].

Nous pouvons donc affirmer qu'aucune méthode n'est parfaite : peu importe la réduction de dimension appliquée une perte de données est inévitable ce qui cause forcément l'apparition d'artefacts. Ce qui soulève la question suivante : Comment s'assurer de la qualité d'une projection ? Pour appréhender cela, nous passerons en revue plusieurs critères et méthodes d'évaluation de la qualité d'une projection dans le chapitre suivant.

## 2 Les critères d'évaluation des projections

Comme nous l'avons vu dans le chapitre précédent, la projection d'un jeu de données multidimensionnel passe par un processus de réduction de dimension. Or, celui-ci a un coût : il implique une perte de données qui va se répercuter lors de la représentation finale de celles-ci. La réduction induit une erreur : le stress qui a un impact sur l'approximation des distances au sein de la représentation.

Par définition, une bonne projection est une projection qui représente le plus fidèlement possible la structure sous-jacente du jeu de données. Par conséquent, une bonne configuration de projection est une configuration qui minimise le stress de la projection. En d'autres termes: plus les distances entre le jeu de données et la projection seront proches, moins il y aura d'erreurs. Ainsi, le but dans une projection est de minimiser le stress.

Mathématiquement parlant, le stress peut être défini comme la somme quadratique des écarts de distances :

$$\epsilon_{\psi} = \sum_{i,j}^n [(d_m(x_i, x_j) - d_p(y_i, y_j))]^2$$

### 2.1 Visualisation du stress

Le stress va produire des artefacts lors de la projection de données. Ceux-ci vont avoir un impact sur l'approximation et l'interprétation de la représentation graphique. Il y a plusieurs façons de le quantifier :

Les mesures peuvent être locales, c'est-à-dire qu'elles sont visualisées en chaque point et aident à expliquer la projection. C'est possible en utilisant un diagramme de Shepard [31] [4]. D'autres méthodes utilisent un *jitter disc* autour de chaque point. Cela permet de visualiser le stress, mais ne donne pas des informations sur les origines de celui-ci [1].

Il est également possible de le visualiser en utilisant un encodage colorimétrique, qui en fonction d'une couleur et de son intensité indiquera les points les moins fiables de la projection [8]. On peut aussi avoir recours à différents types d'interpolation de la couleur entre les points, selon différents paramètres, ce qui forme une carte de précision [32].

Une méthode pour mettre en avant les artefacts topologiques (faux voisinages et déchirures) a été mise en place par Lespinats et Aupetit [8] : l'idée est de visualiser deux mesures de stress simultanément à l'aide d'une échelle de couleur 2D uniforme. Cela permet d'utiliser la projection malgré les erreurs de positionnement des différents points: cependant la mesure de qualité ne garantit pas qu'on puisse exploiter la projection (surtout pour les tâches de clustering visuel) .

Il y a également, des méthodes interactives avec différentes vues possibles (stress) globaux, faux voisinages, déchirures ... selon les différents niveaux de granularité (point seul ou groupes de points) [33] [4]. De plus, en fonction du type d'artefacts, différentes mesures d'erreurs sont proposées : coloration interpolée ou edge bundling (arc compact) pour, par exemple relier les points et faire le parallèle pour montrer les déchirures sur la projection [34] [35] [4]. Les arcs sont aussi colorés de façon différente en fonction de l'intensité des erreurs.

Ce système est interactif, car l'utilisateur peut choisir manuellement un groupe de point à étudier afin de visualiser seulement les artefacts qui lui sont relatifs. Il y a aussi beaucoup de vues statiques qui peuvent l'aider à filtrer.

### 2.2 Les critères de mesure

Pour évaluer la qualité d'une projection, plusieurs approches se basant sur différents critères existant. Ces critères peuvent être relatifs à la structure des données projetées avec le taux d'oublier de la projection (*outlying*), le taux de cluster et à leurs qualités (*visual clustering*, *class consistency*, *cluster separator*, *class density*), et à la forme de la projection (*clumpy*, *skinny*). Nous allons les détailler au cours de ce chapitre.

## Clustering

Le clustering sert à partitionner les données en groupe homogènes. En d'autres termes, cela consiste à évaluer si le jeu de données contient des groupes(ou clusters) ainsi que leur qualité (c'est-à-dire s' ils sont séparés et facilement identifiables). Pour cela, il existe plusieurs critères.

- **Class consistency** : Tout d'abord, il y a la mesure de la *class consistency* [36]. Dans un jeu de données, chaque classe est composée d'un sous-ensemble de l'espace de données qui lui est assigné. Ainsi chaque point est assigné à une classe. L'ensemble des classes du data-space est appelé *class structure*. L'objectif est donc de vérifier si la séparation des classes est correcte en fonction des dimensions étudiées, sachant qu'une bonne représentation visuelle d'une "Class structure" doit être fidèle à cette même structure de classe. Pour déterminer cela, nous pouvons le calculer selon la règle de la *distance to centroid*: c'est-à-dire que chaque point d'une classe doit être à une distance inférieure de son centroïde(i.e le centre de sa classe) qu'à celui d'une autre classe. Or c'est une des propriétés qui est très souvent perdue lors de la projection. Par conséquent, si la *class consistency* est bonne, les classes seront situées à des régions visuellement séparées dans la projection graphique, les clusters seront donc plus facilement discriminables et la projection potentiellement de meilleure qualité.

La class consistency est le plus souvent calculée de la façon suivante :

$$1 - \frac{|\{p \mid \exists j : d(p, \text{centr}(C_k)) \leq d(p, \text{centr}(C_j))\}|}{m}$$

Où  $C_k$  est la classe de  $p$ ,  $\text{centr}(C_k)$  est le centroid de cette même classe.  $m$  le nombre de classe disponible, et  $d(p, \text{centr}(C_k))$  la fonction *distance to centroid*

- **Class density** : Pour justifier de la qualité du clustering d'une projection, nous pouvons également mesurer la densité de classe (ou *class density*). Ce critère évalue la projection en fonction de ses propriété de séparation. Le but est d'identifier les zones où il y a le moins de chevauchement entre les classes.

Pour calculer ce chevauchement entre classes, la méthode passe par l'utilisation d'une représentation continue du jeu de donnée dans laquelle les points appartenants au même cluster forment une image différente. Cela revient à faire en sorte que pour chaque classe il y ai une image distincte avec des densités continues et régulières calculée à partir du voisinage local de chaque point (ou pixel).

Le principe est donc de retenir les resultats avec le moins de chevauchement [37]. Ainsi, lors d'une projection, les points appartenant à un cluster forment une image, c'est -à -dire que chaque classe forme une image. L'algorithme se fait en fonction de la densité basée sur le voisinage : pour chaque pixel, la distance de son voisin le plus proche de la même classe est enregistrée. Puis, la densité locale est calculée dans une sphère de rayon de la distance maximum possible. Le chevauchement global des classes est ensuite estimé en calculant la somme de la différence de chaque paire de pixels à la valeur absolue. Puis la visualisation avec le chevauchement le plus faible sera retenue.

$$CDM = \sum_{k=1}^{m-1} \sum_{l=k+1}^m \sum_{i=1}^P \|P_k^i - P_l^i\|$$

- **Histogram density** : Il y a aussi la mesure de l'*histogram density* [4] [38] qui est une mesure de l'entropie de la projection (i.e l'information moyenne de différentes portions de la projection découpée en grille). La projection est séparée en plusieurs barres (à l'image des histogrammes). Dans chacune de ces barres on comptabilise le nombre de points et leur

étiquette de classe. L'entropie de chaque barre se calcule de cette façon.

$$H(p) = - \sum_c \frac{P_c}{\sum_c P_c} \log_2 \frac{P_c}{\sum_c P_c}$$

Ainsi, si tous les points d'une barre sont de la même classe, alors l'entropie est égale à zéro. Par conséquent le clustering sera de bonne qualité.

## Outlying

Le taux d'outliers est important à définir pour justifier de la qualité du clustering. Un *outlier* est une donnée aberrante. Pour déterminer si un point est un outlier, il y a plusieurs méthodes.

### Critères utilisant les graphes :

- Une des façons de raisonner en graphe et d'utiliser le concept de *Minimum Spanning Tree* (MST) [39], c'est-à-dire l'arbre le moins long que l'on peut créer à partir du nuage de points. En prenant comme référentiel le MST, les outliers sont soit des points qui sont localisés aux extrémités de l'arbre ou en intérieur dans des régions relativement vides. Ainsi, on peut les définir comme des nœuds de degrés 1, et dont la somme des poids des arêtes adjacentes est supérieure à 'w' selon la formule suivante

$$\omega = q_{75} + 1.5(q_{75} - q_{25})$$

. Avec  $q_{75}$  représentant le 75ème percentile de la longueur des arêtes du MST et la partie entre parenthèse représentant l'intervalle interquartile de la longueur des arêtes. Puis, nous pouvons calculer le taux d'outliers dans la projection (outlying) en divisant la longueur totale de tous les outliers par la longueur totale de l'arbre selon la formule suivante :

$$c_{outlying} = length(T_{outliers}) / length(T)$$

- Il y a également la méthode de l'*isolation Forest*. Dans celle-ci, l'algorithme choisit une caractéristique du jeu de données et une valeur *split* comprise entre le maximum et le minimum des valeurs. Cette étape est effectuée pour toutes les observations. Par la suite la moyenne des tous les arbres est faite pour construire la forêt. L'apprentissage de l'algorithme consiste à comparer les observations avec la *splitting value* dans un nœud. Ce nœud comprendra également deux sous-nœuds dans lesquels on fera aussi la même comparaison. Ainsi le nombre de *splitting* est égal à la longueur du chemin. Toutes les comparaisons auront un score allant de 0 à 1. 0 signifiant "normalité", et 1 signifiant qu'il y a un grand nombre d'outlier [40]. Elle est implémentée en python avec scikit Learn.

Puis, nous pouvons calculer le taux d'outliers dans la projection (*outlying*) en divisant la longueur totale de tous les outliers par la longueur totale de l'arbre selon la formule suivante.

$$c_{outlying} = length(T_{outliers}) / length(T)$$

**Le Local Outlier Factor (LOF):** Pour ce critère, on considère un outlier selon son voisinage local [41]. Le LRD (*Local Reachability Density*) de chaque point est comparé avec le LRD de ses voisins. Puis on calcule le ratio de la moyenne des LRD des voisins du point sur le LRD du point lui-même. Si le Local Outlier Factor est supérieur à 1 alors le point est un outlier. Cette mesure est très utile pour détecter des outliers dans les clusters extrêmement denses. Néanmoins cette

valeur étant un ratio et non un seuil, elle est parfois compliquée à interpréter en fonction de la problématique de recherche.

### Critères de formes

Les informations sur la forme de la projection sont également très importantes.[39] La forme est souvent définie par le critère *skinny*. Ce critère représente grossièrement la finesse de la courbe. Celui-ci est défini en calculant le ratio aire/périmètre d'une projection (avec normalisation)

$$C_{skinny} = 1 - \sqrt{4\pi area(A)/perimeter(A)}$$

Dans le cas où le chemin le plus court d'un Minimum Spanning Tree est quasiment aussi long que la somme des arêtes de l'arbre d'origine, la courbe peut être qualifiée de *stringy*

$$C_{stringy} = diameter(T)/length(T)$$

### Critère de densité

Dans une configuration éparpillée (*scattered*) il y a également le critère de densité qui importe. C'est la distribution des arêtes du Minimum Spanning Tree qui va donner des informations sur la densité relative de points. Le calcul de cette densité se base sur la mesure des quantiles de la longueur des arêtes :

$$C_{skew} = (q_{90} - q_{50})/(q_{90} - q_{10})$$

Ainsi, si la dimension d'un *Minimum Spanning Tree* est extrêmement asymétrique, les clusters y seront mal représentés. Dans ce cas, elle est qualifiée de *clumpy*. [39]

### Score de Projection (*seulement pour l'ACP*)

Les informations d'une projection par ACP sont mesurées en comparant la variance totale du jeu de données et celle qui a été retenue après la réduction de dimension. Ce qui la rend très dépendante du nombre de dimensions du jeu de données de base.

Ainsi, sauf pour les jeux de données de petite taille, il est quasiment impossible que les trois composantes principales d'un grand jeu de données retiennent 100 pourcents(ou quasiment 100 pourcents) de la variance. Et ce, même si le jeu de données est composé de variables non aléatoires avec des structures très facilement interprétables. Le but de cette méthode est de quantifier l'informativité de la projection non pas par la variance finale, mais par l'excès de variance de ce qui est attendu après réduction d'un dataset aléatoire. Pour calculer le score de projection: on calcule la variance totale qui est retenue par les 3 principaux composants. Ensuite, on estime la valeur de ces mêmes entités, mais pour un jeu de variables complètement aléatoire.

Enfin, on calcule le score de projection en faisant différence entre la racine carrée de la quantité observée et de la racine carrée de la valeur attendue (celle du random dataset). Ainsi, si le résultat est grand(s'éloigne positivement de zéro) cela veut dire qu'il y a plus d'informations (relatives au calcul de la variance) pour le vrai jeu de données que pour celui du jeu de données aléatoire. Par conséquent, la projection est plus susceptible de proposer des structures fidèles et intéressantes, qui ne sont pas dues au hasard[42].

## 2.3 Conclusion de la partie

Ces mesures de qualité sont donc utiles pour les personnes non spécialistes. Elles leur permettent d'appréhender au mieux la visualisation obtenue après projection. Elles peuvent également être utilisées de façon automatique dans des algorithmes qui permettent de trier et de déterminer si la projection qui va être affichée est fidèle au jeu de données[43].

Comme nous l'avons brièvement abordé dans l'état de l'art, les facteurs numériques à eux seuls ne

suffisent pas à l'interprétation d'une projection : il faut rajouter le facteur humain à l'équation. Malgré l'utilité de toutes ces mesures, le jugement humain reste toujours le plus important[37]. De plus, c'est également l'être humain qui a pour rôle d'interpréter la projection. C'est pour cela qu'il existe des évaluations à faire passer aux humains, qui servent à mesurer l'impact des techniques de projection, la qualité des graphiques et leur compréhension par l'utilisateur. Par exemple ,nous pouvons tester directement les utilisateurs avec plusieurs types de tâches telles que le *data outlier validation*, le *clustering validation*, le *cluster énumération*, et le *class outliers validation* [1]. Ces tests consistent en la détection des structures et des données aberrantes dans les projections. Ainsi si une projection a une bonne qualité (suivant les critères numériques) et permet aux utilisateurs de discriminer rapidement et correctement les structures, elle tend alors vers ce que l'on peut considérer comme étant une bonne projection.



## Réalisation

Comme mentionné dans l'état de l'art la visualisation des grands jeux de donnée souffre de cette *malédiction de la dimensionalité*. En effet, le problème des projections n'est pas celui de projeter les données, mais de les projeter fidèlement. C'est-à-dire, la difficulté pour la méthode à respecter le rapport des distances des points entre le jeu de donnée original (à  $n$ -dimensions) et la projection finale (qui réduit ce dernier à deux ou 3 dimensions). Plus le nombre de dimensions augmente, plus il devient compliqué de respecter les distances après réduction.

Le choix d'une méthode de projection par rapport à une autre dépend d'autant plus de la sémantique et de la structure des données. Ce choix dépendra aussi de ce que l'on souhaite observer : voulons-nous mettre en valeur les classes ? Voulons-nous détecter des valeurs aberrantes ? En fonction de cela, le choix de la méthode et de son paramétrage peut grandement différer. Une fois le jeu de données projeté, il faut s'assurer de la qualité de cette projection. Pour cela nous disposons de nombreux critères qui peuvent donner un indice sur le clustering, outlying et la densité de la projection. Certains critères sont même propres à certaines méthodes (comme le projection score pour l'ACP). Encore une fois, en fonction de ce que nous voulons observer certains critères seront plus utiles que d'autres. De plus, ces critères ne font pas tout : d'une certaine façon, une projection se doit d'être à la fois la plus compréhensible par l'Homme et la plus fidèle possible au jeu de données qu'elle représente.

C'est dans ce but, que dans la partie réalisation de cette UE, nous avons proposé une méthode d'évaluation des différentes méthodes de projections citées. Le but de cette démarche est de proposer un outil capable de comparer sa méthode de projection à d'autres via un tableau de scoring pour se rendre compte de la pertinence et de la fidélité des différentes méthodes de projection. Dans cette partie nous détaillerons l'implémentation, les jeux de données utilisés et le fonctionnement du programme.

## 3 Implémentation

### 3.1 Méthodologie

Nous avons donc mis en place un programme avec un système de score qui classe l'efficacité des différentes méthodes de projection pour un même jeu de données. C'est-à-dire que le jeu de donnée va être projeté une fois par chaque méthode, puis le résultat sera évalué. Pour mesurer ces résultats nous avons utilisé différents critères issus de la librairie R *clusterCrit* [44] et *scikit-learn* [45]. ClusterCrit est une librairie R qui fournit une liste de critères permettant d'attester de la qualité interne des clusters et une liste de critères qui permet de mesurer la similarité entre deux partitions. Ces critères prennent seulement en compte la répartition des points dans les différents cluster et ne permettent pas de mesurer la qualité de la distribution [46]

### 3.2 Jeux de données utilisés

Pour comparer les différentes méthodes nous avons générés plusieurs jeux de données qui présentent les caractéristiques suivantes :

- Un jeu de données avec deux clusters superposés avec un seul (Figure 5)
- Un jeu de données avec trois clusters bien distincts (Figure 6)
- Un jeu de données avec trois clusters superposés (Figure 7)

Pour les créer nous avons utilisé la méthode *makeblobs* de la librairie *sklearn*.

Les paramètres fixes sont :

- *n samples* qui représente le nombre de points total qui sera également réparti entre les différents clusters. Nous l'avons fixé à vingt-mille.
- *n features* qui représente le nombre de variables pour chaque échantillon. Nous l'avons fixé à vingt-et-un.

Nous avons fait varier les paramètres suivants pour obtenir la répartition souhaitée :

- *centers* qui détermine le nombre de centres à générer
- *cluster std* qui détermine l'écart type des clusters
- *random state* qui détermine un nombre aléatoire pour la création du jeu de données.

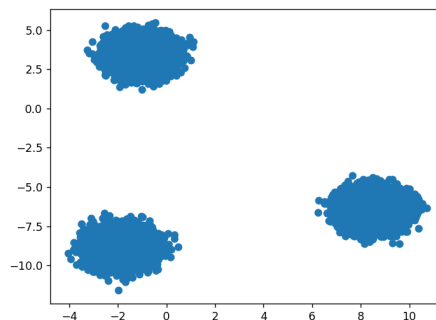


Figure 5: Trois clusters distincts

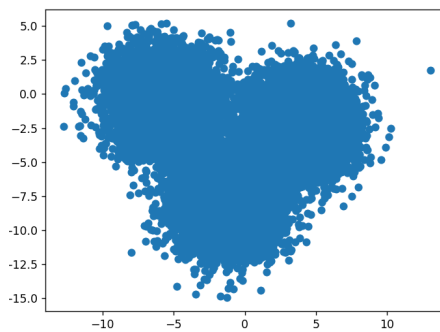


Figure 6: Trois clusters avec chevauchement

### 3.3 Fonctionnement du programme

La librairie *rpy2* est utilisée pour pouvoir charger et utiliser des librairies R sous Python. Comme précisé plus haut nous appliquons cela à la librairie *clusterCrit*. Et appelons tous les critères d'intérêts qui sont listés dans la variable globale *crit*.

La fonction *calculate* va utiliser chaque méthode de réduction sur le jeu de donnée, le resultat de cette réduction sera stockée dans la variable *p2* à laquelle on appliquera la fonction *IntVector*

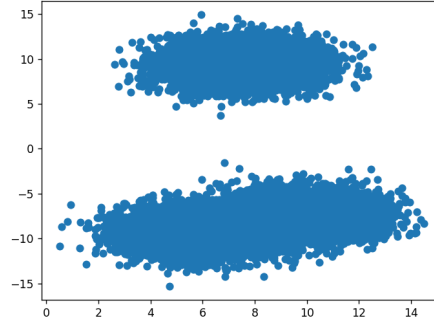


Figure 7: Deux clusters avec chevauchement et un cluster seul

qui permet de convertir un tableau python en vecteur utilisable sous R. Comme éléments de comparaison, on applique Kmean sur le jeu de donnée originel. ce qui permet de partitionner le jeu de données en  $k$  groupes. Dans notre cas  $k = 3$ .

Une fois ces deux opérations effectuées, il y a deux variables *rds* (qui représente le résultat de la méthode de réduction), et *rdsOriginal* (qui représente le jeu de donnée original après kmean.) Pour comparer, ces deux jeux de données nous les passons en paramètre de la fonction *metricsCalcul*. Celle-ci va itérer à travers tous les critères de la liste *crit* et appliquer l'évaluation de chaque critères.

### 3.4 Expérimentation

Les résultats seront stockés sous la forme d'un tableau comprenant les résultats de chaque méthode pour chaque jeux de données. Voici un tableau contenant la valeur minimum et Maximum de chaque critères.

|     | Czekanowski | Folkes | Hubert | Jaccard | Kulczynski | McNemar | Phi      | Precision | Rand | Recall | Rogers Tanimoto | Russel Rao | SokalS1 | SokalS2 |
|-----|-------------|--------|--------|---------|------------|---------|----------|-----------|------|--------|-----------------|------------|---------|---------|
| Min | 0           | 0      | -1     | 0       | 0          | 0       | -1       | 0         | 0    | 0      | 0               | 0          | 0       | 0       |
| Max | 1           | 1      | 1      | 1       | 1          | 1       | $\infty$ | 1         | 1    | 1      | 1               | 1          | 1       | 1       |

Figure 8: Résultats des jeux de données pour l'indice de Hubert

## 4 Analyse des résultats

|               | Czekanowski | Folkes      | Hubert | Jaccard     | Kulczynski  | McNemar       | Phi  | Precision   | Rand        | Recall      | Rogers Tanimoto | Russel Rao  | SokalS1     | SokalS2    | Total score |
|---------------|-------------|-------------|--------|-------------|-------------|---------------|------|-------------|-------------|-------------|-----------------|-------------|-------------|------------|-------------|
| MDS dataset1  | 0.12        | 0.16        | -0.0   | 0.06        | 0.2         | 306.01        | 0.0  | 0.33        | 0.64        | 0.07        | 0.47            | 0.02        | 0.03        | 0.78       | <b>0/14</b> |
| tSNE dataset1 | 0.02        | 0.06        | 0.0    | 0.01        | 0.17        | <b>391.03</b> | 0.0  | <b>0.34</b> | <b>0.66</b> | 0.01        | <b>0.5</b>      | 0.0         | 0.01        | <b>0.8</b> | <b>5/14</b> |
| pca dataset1  | <b>0.31</b> | <b>0.31</b> | -0.0   | <b>0.18</b> | <b>0.31</b> | 55.63         | -0.0 | 0.33        | 0.57        | <b>0.28</b> | 0.4             | <b>0.1</b>  | <b>0.1</b>  | 0.73       | <b>7/14</b> |
| mds dataset2  | 0.08        | 0.12        | -0.0   | 0.04        | 0.19        | 344.52        | 0.0  | 0.33        | 0.65        | 0.04        | 0.48            | 0.01        | 0.02        | 0.79       | <b>0/14</b> |
| tSNE dataset2 | 0.02        | 0.06        | 0.0    | 0.01        | 0.17        | <b>393.52</b> | 0.0  | <b>0.34</b> | <b>0.66</b> | 0.01        | <b>0.5</b>      | 0.0         | 0.0         | <b>0.8</b> | <b>5/14</b> |
| pca dataset2  | <b>0.13</b> | <b>0.17</b> | 0.0    | <b>0.07</b> | <b>0.21</b> | 293.66        | 0.0  | 0.33        | 0.64        | <b>0.08</b> | 0.47            | <b>0.03</b> | <b>0.04</b> | 0.78       | <b>7/14</b> |
| mds dataset3  | 0.04        | 0.08        | -0.0   | 0.02        | <b>0.18</b> | 377.4         | 0.0  | 0.33        | <b>0.66</b> | 0.02        | 0.49            | 0.01        | 0.01        | <b>0.8</b> | <b>3/14</b> |
| tSNE dataset3 | 0.02        | 0.06        | 0.0    | 0.01        | <b>0.18</b> | <b>391.69</b> | 0.0  | <b>0.35</b> | <b>0.66</b> | 0.01        | <b>0.5</b>      | 0.0         | 0.01        | <b>0.8</b> | <b>5/14</b> |
| pca dataset3  | <b>0.06</b> | <b>0.1</b>  | -0.0   | <b>0.03</b> | <b>0.18</b> | 362.47        | -0.0 | 0.33        | <b>0.66</b> | <b>0.03</b> | 0.49            | 0.01        | <b>0.02</b> | 0.79       | <b>7/14</b> |

Figure 9: Résultats des jeux de données pour l'indice de Hubert

### 4.1 Analyse de chaque critères

Dans cette section, nous allons détailler les résultats des jeux de données pour les différentes indices.

### L'index de Czekanowski Dice

On peut voir que pour les trois jeux de données c'est l'ACP qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve de meilleur résultat.

### L'index de Folkes Mallow

On peut voir que pour les trois jeux de données c'est l'ACP qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve de meilleur résultat.

### L'indice de Hubert

Pour cet indice les résultats précis sont les suivants :

|               |           |
|---------------|-----------|
| MDS dataset1  | -0.000366 |
| tSNE dataset1 | 0.000334  |
| PCA dataset1  | 0.000119  |
| MDS dataset2  | -0.001389 |
| tSNE dataset2 | 0.000749  |
| PCA dataset2  | 0.00416   |
| MDS dataset3  | -0.000556 |
| tSNE dataset3 | 0.001285  |
| PCA dataset3  | -0.000369 |

Figure 10: Résultats des jeux de données pour l'indice de Hubert

Nous pouvons voir que les valeurs 0.0 et  $-0.0$  viennent de l'arrondi des résultats. Les résultats à cet indice sont négligeables la corrélation est donc proche de zéro.

### L'indice de Jaccard

On peut voir que pour les trois jeux de données c'est l'ACP qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve de meilleur résultat.

### L'indice de Kulczynski

On peut voir que pour les trois jeux de données c'est l'ACP qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve de meilleur résultat. On peut noter que pour le troisième jeu de données (trois clusters se chevauchant) toutes les méthodes de réduction ont la même efficacité : elles ont toutes les trois eu un score de 0.18.

### L'indice de McNemar

On peut voir que pour les trois jeux de données c'est le tSNE qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve de meilleur résultat. Néanmoins, son score(394.03) est très proche du second jeux de données(393.52).

## L'indice de Phi

Pour cet indice les résultats précis sont les suivants :

|               |               |
|---------------|---------------|
| MDS dataset1  | 9.172429e-12  |
| tSNE dataset1 | 3.508708e-12  |
| PCA dataset1  | -2.322457e-12 |
| MDS dataset2  | -0.001389     |
| tSNE dataset2 | 6.570988e-12  |
| PCA dataset2  | 4.064610e-12  |
| MDS dataset3  | -3.519841e-12 |
| tSNE dataset3 | 1.093924e-11  |
| PCA dataset3  | -2.284271e-12 |

Figure 11: Résultats des jeux de données pour l'indice Phi

Nous pouvons voir que les valeurs 0.0 et  $-0.0$  viennent de l'arrondi des résultats. Les résultats à cet indice ont pour ordre de grandeur  $10^{-12}$ , ils sont donc négligeables.

## L'indice de Jaccard

On peut voir que pour les trois jeux de données c'est l'ACP qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve de meilleur résultat.

## L'indice Précision

On peut voir que pour cet indice c'est le tSNE qui a les meilleurs résultats pour les trois jeux de données. Cependant la différence reste faible par rapport aux autre méthodes (seulement 0.01).

## L'indice Rand

Pour le premier et le second jeu de données le tSNE a un meilleur score que les deux autres techniques. Pour le troisième jeu de données, toutes les techniques semblent se valoir : elles ont tous le même score de 0.66.

## L'indice Recall

Pour cet indice c'est l'ACP qui obtient de meilleurs scores pour les trois jeux de données. On peut noter que la différence de score entre les méthodes est beaucoup plus marquée pour le premier jeu de données.

## L'indice Rogers Tanimoto

Pour cet indice c'est principalement toutes les techniques ont sensiblement le même score. Néanmoins on peut noter une légère tendance pour la tSNE à être plus efficace. Elle est de 0.5 pour les trois jeux de données.

## L'indice Russel Rao

Pour cet indice, nous pouvons voir que pour les trois jeux de données c'est l'ACP qui a les meilleurs résultats pour les trois jeux de données. Parmi les trois jeux de données c'est le premier (deux clusters se chevauchant et un seul) avec lequel on retrouve une différence plus marquée entre les scores des méthodes : le score de l'ACP est à 0.1 tandis que le score du mds est à 0.02 et celui du tSNE est à 0.0.

## Les indices Sokal Sneath

Pour le premier indice de Sokal Sneath nous pouvons voir que c'est l'ACP qui a de meilleurs scores pour tous les jeux de données, et que parmi ceux-ci c'est avec le premier que l'on trouve la plus grande différence de score entre méthode : le score de la PCA est à 0.1 tandis que le score du mds est à 0.03 et celui du tSNE est à 0.1.

Pour le second indice de Sokal Sneath c'est le tSNE qui semble avoir le meilleur score pour les trois jeux de données. On peut noter que pour le troisième jeu de données, les scores du MDS, du tSNE sont les mêmes 0.8 et que celle de l'ACP s'en rapproche de très près étant donné que sa valeur est de 0.79.

## 4.2 Discussion des résultats

Nous pouvons constater que dans une grande majorité des cas c'est l'ACP qui a les meilleurs résultats. De plus, c'est systématiquement le premier jeu de données (deux clusters qui se chevauchent et un cluster seul) pour lequel cette méthode fonctionne le mieux. Néanmoins ces résultats ne permettent pas d'affirmer que l'ACP est la méthode de projection la plus efficace, et ce, pour plusieurs raisons. D'une part il serait judicieux d'ajouter à ce programme d'autres techniques de projection, ce qui permet d'avoir plus d'éléments de comparaison. Tout d'abord, il aurait été intéressant de tester la méthode classNerv (qu'il a été compliqué d'implémenter suite à des problèmes de compatibilité), qui au vu des résultats présentés dans la première partie de ce rapport semble être particulièrement performante et aurait pu montrer des résultats plus probants que ceux de l'ACP.

De plus, il faudrait diversifier les différents jeux de données c'est à dire faire des tests sur d'autres jeux de données générés, par exemple un jeu de données avec trois clusters bien distincts et des outliers, ou bien un jeu de données avec un cluster et plusieurs outliers. Il faudrait également tester le programme avec des jeux de données qui proviennent de sources réelles (par exemple trouvables sur Kaggle ou Datahub.io) et tester différents types de jeux de données (bivariés, multivariés, catégoriels...)

## 5 Conclusion

Pour conclure, nous avons vu que dans une certaine mesure l'ACP se révèle être la méthode la plus efficace pour trois types de jeux de données générés différents. À savoir :

- Un jeu de données avec deux clusters superposés avec un seul (Figure 5)
- Un jeu de données avec trois clusters bien distincts (Figure 6)
- Un jeu de données avec trois clusters superposés (Figure 7)

Il serait intéressant d'utiliser ClassNerv qui en théorie est la plus performante. Nous pourrions par exemple effectuer un étalonnage des différentes autres méthodes de projection, en l'utilisant en tant que référentiel. C'est à dire que lorsqu'un jeu de données serait projeté avec une autre méthode de projection, le résultat obtenu serait comparé à la projection de ce même jeu de données obtenue via ClassNerv. Puis, en fonction des divergences entre ces deux jeux de données projetés, certains biais seraient mesurés et affichés. De plus dans l'objectif d'être le plus précis possible, nous envisageons également de mesurer les différents critères de la projection (clustering, outlying,

forme...)et de croiser les résultats obtenus à ceux de la comparaison des deux projections. Nous pensons aussi qu'il serait judicieux de mettre en place un système de "scoring", qui permettrait de rendre plus intuitive l'appréhension des résultats relatifs à la qualité de la visualisation.

# Bibliography

- [1] N. Heulot, J.-D. Fekete, and M. Aupetit, “Visualizing dimensionality reduction artifacts: An evaluation,” *arXiv preprint arXiv:1705.05283*, 2017.
- [2] M. Card, *Readings in information visualization: using vision to think*. Morgan Kaufmann, 1999.
- [3] D. Keim, G. Andrienko, J.-D. Fekete, C. Görg, J. Kohlhammer, and G. Melançon, “Visual analytics: Definition, process, and challenges,” in *Information visualization*. Springer, 2008, pp. 154–175.
- [4] N. Heulot, “Etude des projections de données comme support interactif de l’analyse visuelle de la structure de données de grande dimension,” Ph.D. dissertation, Université Paris-Sud, 2014.
- [5] K. Koffka, “Principle of gestalt psychology, translated by w. li,” 1997.
- [6] S. Johansson and J. Johansson, “Interactive dimensionality reduction through user-defined combinations of quality metrics,” *IEEE transactions on visualization and computer graphics*, vol. 15, no. 6, pp. 993–1000, 2009.
- [7] M. Aupetit, “Visualizing distortions and recovering topology in continuous projection techniques,” *Neurocomputing*, vol. 70, no. 7-9, pp. 1304–1330, 2007.
- [8] S. Lespinats and M. Aupetit, “Checkviz: Sanity check and topological clues for linear and non-linear mappings,” in *Computer Graphics Forum*, vol. 30, no. 1. Wiley Online Library, 2011, pp. 113–125.
- [9] P. V. Tran and T. X. Le, “Approaching human vision perception to designing visual graph in data visualization,” *Concurrency and Computation: Practice and Experience*, vol. 33, no. 2, p. e5722, 2021.
- [10] M. F. Deering, “The limits of human vision,” in *2nd International Immersive Projection Technology Workshop*, vol. 2, 1998, p. 1.
- [11] G. Dzemyda, V. Medvedev, A. Lupeikiene, O. Kurasova, and A. Caplinskas, “Big multidimensional datasets visualization using neural networks—efficient decision support,” *Complex Systems Informatics and Modeling Quarterly*, no. 6, pp. 1–11, 2016.
- [12] R. L. Somorjai, B. Dolenko, A. Nikulin, W. Roberson, and N. Thiessen, “Class proximity measures—dissimilarity-based classification and display of high-dimensional data,” *Journal of biomedical informatics*, vol. 44, no. 5, pp. 775–788, 2011.
- [13] A.-H. Karimi, A. Wong, and A. Ghodsi, “Srp: Efficient class-aware embedding learning for large-scale data via supervised random projections,” *arXiv preprint arXiv:1811.03166*, 2018.
- [14] J.-N. Hwang, H. Li, M. Maechler, R. D. Martin, and J. Schimert, “A comparison of projection pursuit and neural network regression modeling,” in *NIPS*, 1991, pp. 1159–1166.



- [15] T. Dickhaus, “Simultaneous statistical inference,” in *With applications in the life sciences*. Springer, 2014.
- [16] B. Colange, J. Peltonen, M. Aupetit, D. Dutykh, and S. Lespinats, “Steering distortions to preserve classes and neighbors in supervised dimensionality reduction,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 13 214–13 225. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/99607461cdb9c26e2bd5f31b12dcf27a-Paper.pdf>
- [17] H. Abdi, “Metric multidimensional scaling (mds): analyzing distance matrices,” *Encyclopedia of measurement and statistics*, pp. 1–13, 2007.
- [18] G. Hinton and S. T. Roweis, “Stochastic neighbor embedding,” in *NIPS*, vol. 15. Citeseer, 2002, pp. 833–840.
- [19] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne.” *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [20] G. C. Linderman and S. Steinerberger, “Clustering with t-sne, provably,” *SIAM Journal on Mathematics of Data Science*, vol. 1, no. 2, pp. 313–332, 2019.
- [21] C. Chui and D. Donoho, “Special issue on diffusion maps and wavelets,” pp. 1–2, 2006.
- [22] U. Von Luxburg, “A tutorial on spectral clustering,” *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [23] J. W. Sammon, “A nonlinear mapping for data structure analysis,” *IEEE Transactions on computers*, vol. 100, no. 5, pp. 401–409, 1969.
- [24] P. Demartines and J. Hérault, “Curvilinear component analysis: A self-organizing neural network for nonlinear mapping of data sets,” *IEEE Transactions on neural networks*, vol. 8, no. 1, pp. 148–154, 1997.
- [25] K. Q. Weinberger and L. K. Saul, “An introduction to nonlinear dimensionality reduction by maximum variance unfolding,” in *AAAI*, vol. 6, 2006, pp. 1683–1686.
- [26] L. Song, A. J. Smola, K. M. Borgwardt, A. Gretton *et al.*, “Colored maximum variance unfolding,” in *Nips*. Citeseer, 2007, pp. 1385–1392.
- [27] L. K. Saul and S. T. Roweis, “An introduction to locally linear embedding,” *unpublished*. Available at: <http://www.cs.toronto.edu/~roweis/lle/publications.html>, 2000.
- [28] M. Belkin and P. Niyogi, “Semi-supervised learning on manifolds,” *Machine Learning Journal*, vol. 1, 2002.
- [29] L. G. Nonato and M. Aupetit, “Multidimensional projection for visual analytics: Linking techniques with distortions, tasks, and layout enrichment,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 8, pp. 2650–2673, 2018.
- [30] F. J. García Fernández, M. Verleysen, J. A. Lee, I. Díaz Blanco *et al.*, “Stability comparison of dimensionality reduction techniques attending to data and parameter variations,” in *Eurographics Conference on Visualization (EuroVis)(2013)*. The Eurographics Association, 2013.
- [31] J. B. Kruskal, “Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis,” *Psychometrika*, vol. 29, no. 1, pp. 1–27, 1964.

How to evaluate the faithfulness of visual data projection ?

- [32] T. Schreck, T. Von Landesberger, and S. Bremm, "Techniques for precision-based visual analysis of projected data," *Information Visualization*, vol. 9, no. 3, pp. 181–193, 2010.
- [33] M. M. Abbas, M. Aupetit, M. Sedlmair, and H. Bensmail, "Clustme: A visual quality measure for ranking monochrome scatterplots based on cluster patterns," in *Computer Graphics Forum*, vol. 38, no. 3. Wiley Online Library, 2019, pp. 225–236.
- [34] D. Holten, "Hierarchical edge bundles: Visualization of adjacency relations in hierarchical data," *IEEE Transactions on visualization and computer graphics*, vol. 12, no. 5, pp. 741–748, 2006.
- [35] D. Holten, P. Isenberg, J. J. Van Wijk, and J.-D. Fekete, "An extended evaluation of the readability of tapered, animated, and textured directed-edge representations in node-link graphs," in *2011 IEEE Pacific Visualization Symposium*. IEEE, 2011, pp. 195–202.
- [36] E. Bertini, A. Tatu, and D. Keim, "Quality metrics in high-dimensional data visualization: An overview and systematization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 12, pp. 2203–2212, 2011.
- [37] A. Tatu, P. Bak, E. Bertini, D. Keim, and J. Schneidewind, "Visual quality metrics and human perception: an initial study on 2d projections of large multidimensional data," in *Proceedings of the International Conference on Advanced Visual Interfaces*, 2010, pp. 49–56.
- [38] A. Tatu, G. Albuquerque, M. Eisemann, J. Schneidewind, H. Theisel, M. Magnork, and D. Keim, "Combining automated analysis and visualization techniques for effective exploration of high-dimensional data," in *2009 IEEE Symposium on Visual Analytics Science and Technology*. IEEE, 2009, pp. 59–66.
- [39] L. Wilkinson, A. Anand, and R. Grossman, "Graph-theoretic scagnostics," in *Information Visualization, IEEE Symposium on*. IEEE Computer Society, 2005, pp. 21–21.
- [40] M. U. Togbe, M. Barry, A. Boly, Y. Chabchoub, R. Chiky, J. Montiel, and V.-T. Tran, "Anomaly detection for data streams based on isolation forest using scikit-multiflow," in *International Conference on Computational Science and Its Applications*. Springer, 2020, pp. 15–30.
- [41] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, "Lof: identifying density-based local outliers," in *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, 2000, pp. 93–104.
- [42] M. Fontes and C. Soneson, "The projection score-an evaluation criterion for variable subset selection in pca visualization," *BMC bioinformatics*, vol. 12, no. 1, pp. 1–17, 2011.
- [43] E. Bertini, A. Tatu, and D. Keim, "Quality metrics in high-dimensional data visualization: An overview and systematization," *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 12, pp. 2203–2212, 2011.
- [44] B. Desgraupes and M. B. Desgraupes, "Package 'clustercrit'," 2018.
- [45] W. T. SCIKIT-LEARN, "Pdf documentation-scikit-learn."
- [46] B. Desgraupes, "Clustering indices," *University of Paris Ouest-Lab Modal'X*, vol. 1, p. 34, 2013.